



PHD

**Generalized Additive Models for Large datasets: spatial-temporal modelling of the UK's Daily Black Smoke (1961 - 2005)**

Li, Zheyuan

*Award date:*  
2019

*Awarding institution:*  
University of Bath

[Link to publication](#)

**Alternative formats**

If you require this document in an alternative format, please contact:  
[openaccess@bath.ac.uk](mailto:openaccess@bath.ac.uk)

Copyright of this thesis rests with the author. Access is subject to the above licence, if given. If no licence is specified above, original content in this thesis is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International (CC BY-NC-ND 4.0) Licence (<https://creativecommons.org/licenses/by-nc-nd/4.0/>). Any third-party copyright material present remains the property of its respective owner(s) and is licensed under its existing terms.

**Take down policy**

If you consider content within Bath's Research Portal to be in breach of UK law, please contact: [openaccess@bath.ac.uk](mailto:openaccess@bath.ac.uk) with the details. Your claim will be investigated and, where appropriate, the item will be removed from public view as soon as possible.

*Citation for published version:*

Li, Z 2018, 'Generalized Additive Models for Large datasets: spatial-temporal modelling of the UK's Daily Black Smoke (1961 - 2005)', Ph.D., University of Bath.

*Publication date:*

2018

*Document Version*

Publisher's PDF, also known as Version of record

[Link to publication](#)

*Publisher Rights*

Unspecified

## University of Bath

**General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

**Generalized Additive Models**  
**for Large datasets:**  
**spatial-temporal modelling**  
**of the UK's**  
**Daily Black Smoke**  
**(1961 - 2005)**

submitted by

**Zheyuan Li**

for the degree of Doctor of Philosophy

of the

**University of Bath**

Department of Mathematical Sciences

April, 2018

**COPYRIGHT**

Attention is drawn to the fact that copyright of this thesis rests with the author. A copy of this thesis has been supplied on condition that anyone who consults it is understood to recognise that its copyright rests with the author and that they must not copy it or use material from it except as permitted by law or with the consent of the author.

This thesis may be made available for consultation within the University Library and may be photocopied or lent to other libraries for the purposes of consultation with effect from .....

Signature of Author .....

Signed on behalf of the School of Mathematics .....

## Abstract

The UK Black Smoke monitoring network has produced daily particulate air pollution data from a network of up to 1200 monitoring stations over several decades, resulting in 10 million measurements in total. Spatial-temporal modelling of the data is desirable for accurate trend / seasonality estimation and mapping and to provide daily exposure estimates for epidemiological cohort studies. Generalized additive models offer one way to do this if we can deal with the data volume and model size. This thesis will develop computation method for estimating generalized additive models having  $O(10^4)$  coefficients and  $O(10^8)$  observations. The strategy combines 3 elements: (i) fine scale discretization of covariates, (ii) an efficient approach to restricted likelihood optimization, that avoids computation of numerically awkward log determinant terms and (iii) restricted likelihood optimization algorithms that make good use of numerical linear algebra methods with high performance and good parallel scaling on modern multi-core machines. The new method enables us to estimate spatial-temporal models for daily Black Smoke data over the last four decades at a daily resolution which had once been infeasible. A spatial-temporal dataset of daily Black Smoke is also produced on a grid of  $5\text{km} \times 5\text{km}$  resolution. Our prediction is shown to suffer from little extrapolation and bias.

# Contents

<b>Introduction</b>	<b>vi</b>
<b>1 Penalized regression splines and additive models</b>	<b>4</b>
1.1 Smoothing by splines . . . . .	4
1.1.1 Low-rank penalized cubic regression spline . . . . .	4
1.1.2 Low-rank penalized thin-plate regression spline . . . . .	7
1.1.3 Low-rank penalized tensor product regression spline . . . . .	8
1.1.4 Canonical form of a low-rank penalized regression spline . . . . .	9
1.1.5 Bayesian interpretation of smoothing . . . . .	10
1.1.6 Smoothness selection via REML . . . . .	11
1.2 Additive models via penalized regression splines . . . . .	12
<b>2 Additive models with splines in practice: simulated examples</b>	<b>13</b>
2.1 Classic example with independent errors . . . . .	13
2.1.1 Tensor product spline for interaction . . . . .	14
2.1.2 Why is an additive model better than a simple smoothing spline? . . . . .	18
2.1.3 Knots placement and choosing number of knots . . . . .	18
2.2 Time series example with AR(1) autocorrelation . . . . .	21
2.2.1 How to do exploratory analysis for data with autocorrelation . . . . .	22
2.2.2 AR(1) error with autocorrelation coefficient $\omega$ . . . . .	22
2.2.3 Golden-section search for point estimation of $\omega$ . . . . .	24

2.3	Visualization of a tensor product spline . . . . .	25
2.4	Model checking . . . . .	27
2.5	Summary . . . . .	28
<b>3</b>	<b>Preliminary modelling of log Black Smoke (logBS)</b>	<b>30</b>
3.1	Introduction to daily logBS dataset and monitoring network . . . . .	31
3.1.1	The daily logBS dataset . . . . .	31
3.1.2	The Black Smoke monitoring network . . . . .	32
3.2	Time series models for daily logBS from “MANCHESTER 11” . . . . .	33
3.2.1	Description of data . . . . .	33
3.2.2	Building a separate model for each day of week . . . . .	35
3.2.3	A joint model for all days of week . . . . .	36
3.2.4	Including meteorological covariates . . . . .	40
3.2.5	Model checking and visualization . . . . .	44
3.2.6	Summary . . . . .	48
3.3	Motivating spatial-temporal modelling . . . . .	52
3.4	Spatial-temporal modelling for yearly mean logBS . . . . .	52
3.4.1	Description of data . . . . .	52
3.4.2	A basic model . . . . .	53
3.4.3	Including i.i.d. random effects . . . . .	56
3.4.4	Model summary and checking . . . . .	57
3.4.5	Visualization of model prediction . . . . .	58
3.4.6	Summary . . . . .	65
3.5	Spatial-temporal modelling for daily logBS in 1967 . . . . .	65
3.5.1	Description of data . . . . .	65
3.5.2	Building a model for Monday data . . . . .	66

3.5.3	Model summary and checking . . . . .	67
3.5.4	Removing spatial autocorrelation in residuals with three-way interaction . . . .	69
3.5.5	Encountering computational hurdle . . . . .	81
3.6	Summary . . . . .	81
<b>4</b>	<b>GAM computation for large datasets: a review</b>	<b>83</b>
4.1	Preliminaries on matrix computations . . . . .	84
4.1.1	Householder QR factorization . . . . .	84
4.1.2	Cholesky factorization . . . . .	86
4.1.3	Solving a triangular system of linear equations . . . . .	87
4.1.4	Positive-definite matrix inverse . . . . .	88
4.1.5	Symmetric eigen decomposition . . . . .	89
4.1.6	Summary . . . . .	90
4.2	Model matrix reduction for additive models . . . . .	93
4.2.1	QR reduction . . . . .	94
4.2.2	Pseudo QR reduction . . . . .	95
4.2.3	Parallel computing . . . . .	96
4.2.4	Summary . . . . .	99
4.3	REML estimation for additive models . . . . .	101
4.3.1	Initial reparametrization . . . . .	103
4.3.2	Derivatives of $\log  \mathbf{S}_{\boldsymbol{\lambda}} _+$ . . . . .	104
4.3.3	Estimation of $\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}$ . . . . .	105
4.3.4	Computing unscaled covariance matrix $\mathbf{V}_{\boldsymbol{\lambda}}$ . . . . .	106
4.3.5	Derivatives of $D_p(\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}, \boldsymbol{\lambda})$ . . . . .	106
4.3.6	Derivatives of $\log  \mathbf{H}_{\boldsymbol{\lambda}} $ . . . . .	108
4.3.7	Special numerical issue: evaluation of $\log  \mathbf{S}_{\boldsymbol{\lambda}} _+$ . . . . .	109

4.3.8	Summary . . . . .	112
4.4	Estimation of GAMs via “performance iteration” . . . . .	112
4.5	Summary . . . . .	114
<b>5</b>	<b>Optimizing GAM computations for high performance computing</b>	<b>115</b>
5.1	A better implementation of the existing computational engine . . . . .	117
5.1.1	How $\log \mathbf{S}_{\lambda,+}$ can be improved . . . . .	117
5.1.2	How $D_p(\hat{\beta}_{\lambda}, \lambda)$ can be improved . . . . .	117
5.1.3	How $\log \mathbf{H}_{\lambda}$ can be improved . . . . .	118
5.2	Using optimized BLAS for high performance computing . . . . .	118
5.2.1	An introduction to BLAS . . . . .	118
5.2.2	Benchmarking BLAS . . . . .	120
5.2.3	Block algorithms and data caching . . . . .	120
5.2.4	Understanding the performance difference between LAPACK routines . . . . .	122
5.2.5	Using optimized BLAS for GAM computations . . . . .	123
5.3	Computing $\mathbf{V}_{\lambda}$ with BLAS . . . . .	123
5.4	Computing $\hat{\beta}_{\lambda}$ via Cholesky factorization . . . . .	124
5.5	Using parallel computing for GAM computations . . . . .	124
5.5.1	Parallel performance of BLAS and LAPACK . . . . .	124
5.5.2	Experimenting parallel computing for GAM computations . . . . .	126
5.6	A comparison of computational capability with INLA . . . . .	127
<b>6</b>	<b>Preliminary modelling of logBS (a revisit)</b>	<b>129</b>
6.1	A revisit to daily logBS model in 1967 . . . . .	129
6.1.1	Justifying three-way interactions with cross-validation . . . . .	129
6.1.2	A joint model for all days of week? . . . . .	129



6.1.3	Block AR(1) errors with autocorrelation coefficient $\omega$ . . . . .	132
6.1.4	Speeding up golden-section search . . . . .	133
6.1.5	Summary and visualization of seven separate models . . . . .	134
6.1.6	Summary . . . . .	137
6.2	A revisit to annual mean logBS model . . . . .	138
<b>7</b>	<b>A model for all daily logBS?</b>	<b>141</b>
7.1	A test model . . . . .	141
7.2	Fitting the test model . . . . .	143
7.3	Checking the test model . . . . .	144
7.4	Summary . . . . .	144
<b>8</b>	<b>Discrete pseudo QR reduction</b>	<b>147</b>
8.1	Faster model matrix formation for pseudo QR reduction . . . . .	147
8.2	Covariate discretization . . . . .	149
8.3	Computation of $\mathbf{X}'\mathbf{W}\mathbf{y}$ . . . . .	149
8.4	Computation of $\mathbf{X}'\mathbf{W}\mathbf{X}$ . . . . .	151
8.5	Computation of $\hat{\mathbf{y}}$ . . . . .	153
8.6	Caching for computation of $\mathbf{X}'\mathbf{W}\mathbf{X}$ . . . . .	155
8.7	Experiment on daily logBS model . . . . .	158
8.8	Summary . . . . .	158
<b>9</b>	<b>Bamboos: boosting discrete pseudo QR reduction by another magnitude</b>	<b>160</b>
9.1	A basic description of bamboos algorithms . . . . .	160
9.1.1	$\mathbf{X}'_i\mathbf{W}\mathbf{X}_j$ between singletons . . . . .	161
9.1.2	$\mathbf{X}'_i\mathbf{W}\mathbf{X}_j$ with tensors . . . . .	163
9.1.3	A unified representation of $\mathbf{X}'_i\mathbf{W}\mathbf{X}_j$ computation . . . . .	164

9.2	FLOP comparison with previous discrete algorithms . . . . .	164
9.2.1	FLOP count for <b>bamboos</b> algorithms . . . . .	165
9.2.2	A recap of FLOP count for the previous discrete algorithms . . . . .	166
9.2.3	FLOP count comparison and summary . . . . .	167
9.3	Fast weights aggregation for $\bar{\mathbf{W}}$ and its sparse construction . . . . .	168
9.3.1	Nested sorting, run-length encoding and bin aggregation . . . . .	168
9.3.2	Sorting algorithms in <b>bamboos</b> . . . . .	170
	Counting sort algorithm for primary sort . . . . .	172
	Binary MSD radix sort for secondary sort . . . . .	173
9.3.3	Dealing with sparse $\bar{\mathbf{W}}$ . . . . .	174
	Compressed column storage . . . . .	174
	Matrix multiplication with sparse $\bar{\mathbf{W}}$ . . . . .	175
	“Dense-sparse switch” in <b>bamboos</b> . . . . .	176
9.3.4	Tri-diagonal $\mathbf{W}$ . . . . .	177
9.4	Experimenting <b>bamboos</b> on daily logBS model . . . . .	178
9.5	Summary . . . . .	178
	<b>Discussion</b>	<b>181</b>
	<b>A Hardware information</b>	<b>184</b>
	<b>Bibliography</b>	<b>191</b>

# Introduction

Since the *Industrial Revolution*, problems associated with air pollution worsened in many countries. During the first half of the 20th century, major pollution episodes occurred in London, notably in 1952 an episode of fog, in which levels of Black Smoke (arising from the incomplete combustion of carbonaceous matter: solid fuels, fuel oils and waste materials, and is especially hazardous to human health in combination with other pollutants which adhere to the particulate surfaces) exceeded  $4500 \mu\text{gm}^{-3}$ , was associated with 4000 excess deaths (Ministry of Health, 1954). Other early episodes, which were caused by a combination of industrial pollution sources and adverse weather conditions, and resulted in large numbers of deaths among the surrounding populations include those in the Meuse valley (Firket, 1936) and the US (Ciocco and Thompson, 1961). Attempts to measure levels of air pollution in a regular and systematic way arose as a result of these episodes.

Following the first *Clean Air Act* in 1952 and a few years to standardize measurement methods (British Standard 1747), the UK established in 1961 the world's first co-ordinated national air pollution monitoring network, the *National Survey of Air Pollution* (Clifton, 1964), to monitor Black Smoke and  $\text{SO}_2$ , with objectives to **1)** monitor progress of Clean Air Acts, and provide the technical basis for future legislation; **2)** monitor progress of Clean Air Acts, and provide the technical basis for future legislation; **3)** provide a consistent body of data for research (For example, effects of climate, topography, industrialization, population density, fuel utilization and urban development can be assessed on pollution levels. Effects of pollution levels on health can also be examined). Since then, the network had worked for 45 years, with a peak size of over 1200 sites per year, producing over 10 million daily measurements.

Because of the size of the dataset, previous statistical modelling with Black Smoke has focussed on modelling time or space averages pollution levels. Fanshawe et al. (2008) studied weekly Black Smoke in Newcastle area (with 20 sites) between 1962 and 1991; Dadvand et al. (2011) studied weekly Black Smoke in North East (a governmental region in England) between 1985 and 1996; Gulliver et al. (2011) produced spatial mapping over Great Britain on a yearly basis between 1962 and 1991; Shaddick and Zidek (2014) studied yearly averaged Black Smoke over 1466 sites across Great Britain between 1966 and 1996. Aggregated air pollution data are not entirely satisfactory from an epidemiological perspective. The variations in pollution are so rapid that even a week is too long a period for true contrasts and comparisons to be appreciated. Besides, acute respiratory disease is usually sensitive to exposure to high levels of pollution over short time periods, and such exposure can be completely hidden in an annual average. Retrospective cohort studies, for example, really require estimates of exposure at the daily level.

There could be several ways to obtain a daily estimate. The simplest one is the running mean, but this only gives estimates for a specific station or location, with no information elsewhere. For spatial prediction purpose, then a spatial smoothing on a daily basis may serve well, but this obscures the variation in time. Particularly, the network itself was changing all the time. There had been nearly 3000 historical stations, but the network only retained a good density over the island for 15 years (1966 - 1980) with about 1000 sites per year. To mitigate these effects, a full spatial-temporal model is most likely to yield reasonable predictions.

At the start of my PhD programme, it was not clear what a model would be like for Black Smoke and what complexity it might have. Preliminary modelling started from modelling small subsets of the dataset, like time series analysis, annual mean, etc. It soon turned out that many effects, like trend, seasonality and temperature effects all vary with space. Building models for these interactions results in great number of parameters. Given such complexity, model representation and data storage are readily huge challenges. Conventional estimation procedure requires “first generating and storing everything” is clearly infeasible. Some estimation routines, like the `bigglm` (Miller, 1992) or `speedglm` (Enea, 2009), even if they do produce smaller fitted model object or a faster estimation, still require storing the whole design matrix for all data and all parameters. In this aspect, the `bam` function from R package `mgcv`, designed for estimating generalized additive models (GAMs) with large datasets, gives a better strategy. Wood et al. (2015) proposed a way to split the huge design matrix into smaller, storable chunks and work with them one at a time. At first glance, this shares the same idea with `bigglm`, but a key feature of `bam` is that these chunks are generated “on the fly”. Once a chunk is processed, it is discarded, making room for the next chunk. The memory footprint is bounded by the size of a chunk, independent of the size of the whole model matrix. For this reason, building GAMs was chosen as the primary method for model development.

This thesis will be structured as follows.

Chapter 1 and Chapter 2 are gentle introductions to smoothing problems and additive models via splines. The former focuses on explaining basic concepts, and the latter demonstrates how additive models can be used in practice via simulated toy datasets.

Chapter 3 conducts some preliminary modelling of Black Smoke by investigating the dataset from different angles: time series, annual mean and daily data from a single year. It turned out that black smoke models can end up with 10000 or even more regression coefficients, making a big computational hurdle for practical model development with GAM. While some statistical questions were also observed during model development, they all gave way to improving GAM fitting methods as without being able to fit a model as fast as possible, it is impractical to do for example, cross-validation.

To see where GAM fitting methods can be improved, it is a prerequisite to understand in depth how existing GAM computational engine works. Chapter 4 is a thorough walkthrough on this.

Chapter 5 starts the first round of performance upgrade of GAM fitting. Many steps are taken to progressively speedup GAM fitting, by the end of which, the computational hurdle previously seen in Chapter 3 are all overcome.

Chapter 6 revisits the preliminary modelling of Black Smoke. Some unanswered questions in Chapter 3 will be addressed.

Chapter 7 makes an attempt at building a daily model for the complete Black Smoke dataset. However, even after using 22000+ parameters, the model is still very inadequate. I eventually give up building models and focus on developing new GAM computation methods for large datasets.

Chapter 8 is a second round of performance upgrading of GAM fitting methods. A novel covariate discretization procedure, packed storage for model matrix, and a set of discrete algorithms to perform required matrix cross-product  $\mathbf{X}'\mathbf{X}$  are developed. A discrete algorithm is actually not as abstract as it sounds. For example, suppose we have a vector of repeated values, like 2, 4, 4, 5, 2, 5, 2 and have compressed it as say a frequency table, where 2, 4, 5 are unique values and 3, 2, 2 are frequency of those unique values. Then computing the mean of the vector does not need to scan through the original vector; it is just a weighted mean of the unique values. Of course, the discrete algorithm used for GAM fitting is not as simple as computing mean; but they share the same idea: store data in a compressed manner and do computation using data in this compressed format. This is not only memory efficient, but also computationally cheaper.

While the discrete algorithm is computationally cheaper, it does not practically deliver high performance. Chapter 9 will develop a new computational engine that is truly high performance computing. With such method, GAM fitting for  $10^7$  or even  $10^8$  data and  $2 \times 10^5$  coefficients can be completed in an hour. I call this new computation methods **bamboos**.

It is indeed a great pity that I am unable to build an adequate daily logBS model. In this way I am missing / bypassing many statistical question that are interesting to many statisticians and worth discussing, like preferential sampling. However, I increasingly realize that building model are not where my passion is. I incline to invest more on **bamboos** to make it a standalone computational module / package that can facilitate, not just computations in **mgcv**, but also general GLM fitting.

# Chapter 1

## Penalized regression splines and additive models

`mgcv` constructs a GAM using low-rank penalized regression splines, which originates from smoothing problems. §1.1 gives an introduction to this, with a focus on the types of splines that are to be applied in this thesis, namely cubic regression splines, thin-plate splines and tensor product splines. Some other aspects of penalized splines are also covered, particularly its link to mixed models and Bayesian inference, which motivates the use of restricted maximum likelihood (REML) for estimating smoothing parameters.

§1.2 moves on to construct a GAM using splines. But since all models in this thesis are additive models, a subclass of GAM, the demonstration will first be restricted to this aspect.

### 1.1 Smoothing by splines

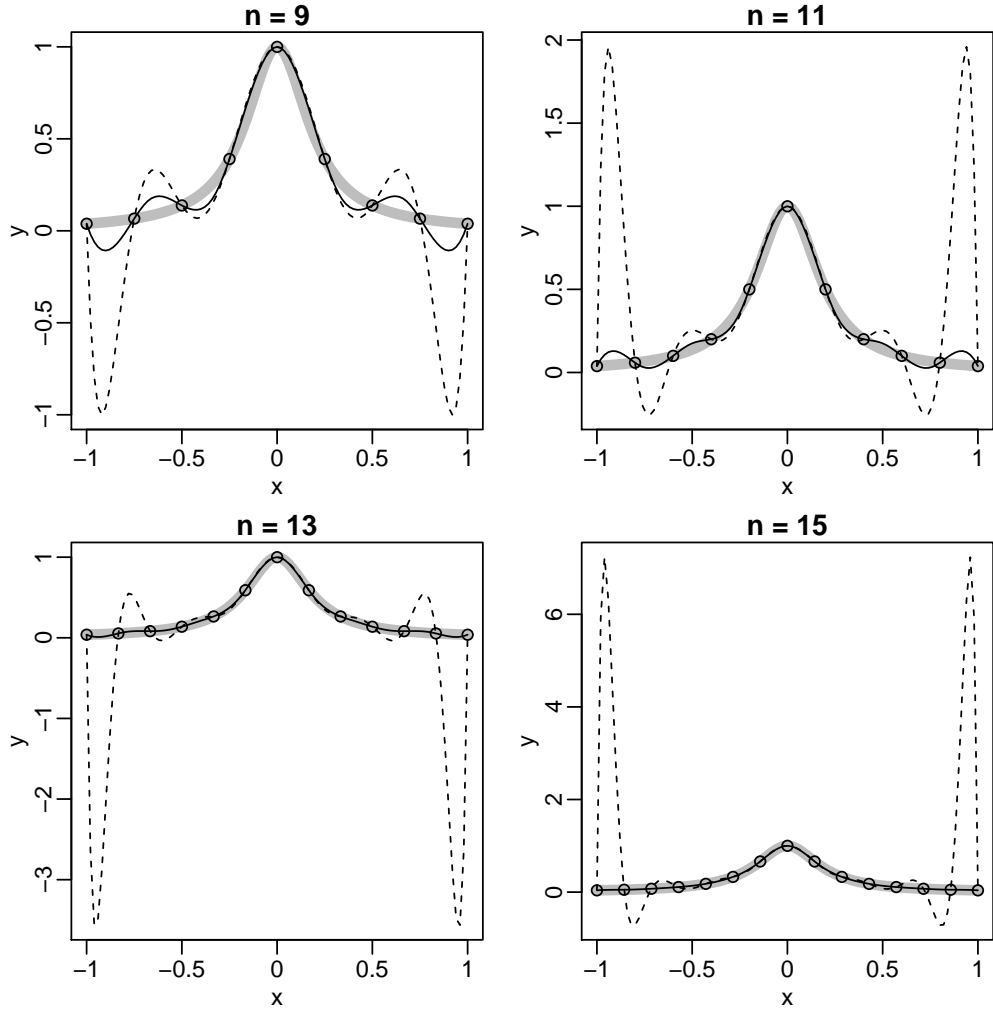
A smoothing problem is a general regression problem where we want to estimate the underlying data generating process from observed data  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  with i.i.d. Gaussian noise. This can be mathematically formulated as

$$y_i = f(x_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \phi/w_i),$$

where ‘ $N$ ’ is Gaussian (or normal) distribution,  $\phi$  is its variance (often denoted by  $\sigma^2$  in literature elsewhere) and  $w_i$  is a known weight for the  $i^{\text{th}}$  datum. Particularly, the shape of  $f(x)$  is not restricted and should be purely “data-driven”. If  $f(x)$  has a pre-specified parametric representation, like a straight line  $f(x) = \beta_0 + \beta_1 x$  or generally a  $p^{\text{th}}$  order polynomial  $f(x) = \sum_{j=0}^p \beta_j x^j$ , then the problem reduces to a linear regression problem (Faraway, 2004). Well-known estimation methods for  $f(x)$  include kernel smoothing (Nadaraya, 1964; Watson, 1964; Wand and Jones, 1994; Bowman and Azzalini, 1997; Hastie et al., 2009), local linear and polynomial regression (Cleveland, 1979; Cleveland and Devlin, 1988; Fan and Gijbels, 1996) and smoothing splines or penalized regression splines to be introduced in this section.

#### 1.1.1 Low-rank penalized cubic regression spline

A spline (de Boor, 1978) is typically a piecewise low-order polynomial that joins smoothly at their connection points, known as *knots*. “Smoothly” means that if the polynomial has order  $p$ , then the  $(p - 1)^{\text{th}}$  derivative of the spline should be continuous at knots. Splines are very useful for



**Figure 1.1:** Illustration of Runge's phenomenon when using Lagrange polynomial of degree  $(n - 1)$  to interpolate the (scaled) Runge function  $R(x) = 1 / (1 + x^2)$ ,  $x \in [-1, 1]$  with  $n$  evenly spaced sample points. In all graphs, the thick gray curve denotes  $R(x)$ , the dashed curve denotes Lagrange interpolation polynomial and the dotted curve is a cubic interpolation spline. A higher order polynomial demonstrates increasingly stronger oscillation at two edges as  $n$  grows. The approximation polynomial does not uniformly converge to the true function, i.e., a higher order polynomial does not improve approximation. By contrast, an interpolation spline (piecewise cubic polynomials) achieves uniform convergence.

interpolation (Kress, 1996), because they do not suffer from the *Runge's phenomenon* (Epperson, 1987) that high-order polynomial interpolation has. See Figure 1.1 for an illustration.

Splines were first introduced to smoothing problems by Reinsch (1967, 1971), who had a variational problem equivalent to

$$\hat{f}(x) = \arg \min_{f(x)} \left\{ \sum_{i=1}^n w_i (y_i - f(x_i))^2 + \lambda \int_{x_1}^{x_n} f''(x)^2 dx \right\},$$

where  $\lambda^{-1}$  is a Lagrangian parameter. Later in statistical literature, the above is termed a *penalized least squares* problem, and  $\int f''(x)^2 dx$  and  $\lambda \geq 0$  are respectively known as a *penalty* of  $f(x)$  and a *smoothing parameter* that controls the strength of the penalization. Reinsch showed that the solution  $f(x)$  is a cubic spline, and then parametrized it by

$$f(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3, \quad x \in [x_i, x_{i+1}),$$

for  $i = 1, 2, \dots, n$  (assuming  $x_{n+1} = +\infty$ ). Under continuity requirement of  $f(x)$ ,  $f'(x)$  and  $f''(x)$

at all  $x_i$ , it can be shown that these parameters are subject to the following constraints:

$$d_i = \frac{h_i}{3}(c_{i+1} - c_i), \quad b_i = \frac{a_{i+1} - a_i}{h_i} - \left(\frac{2}{3}c_i + \frac{1}{3}c_{i+1}\right)h_i, \\ \frac{1}{3}h_i c_i + \frac{2}{3}(h_i + h_{i+1})c_{i+1} + \frac{1}{3}h_{i+1}c_{i+2} = \frac{1}{h_i}a_i - \left(\frac{1}{h_i} + \frac{1}{h_{i+1}}\right)a_{i+1} + \frac{1}{h_{i+1}}a_{i+2},$$

where  $h_i = x_{i+1} - x_i$ . In above, there are  $4n$  parameters for  $n$  data under  $3(n-1)$  constraints. Additional boundary constraints, like *natural boundary condition*  $c_1 = c_n = d_n = 0$  and *periodic boundary condition*  $a_1 = a_n$ ,  $b_1 = b_n$ ,  $c_1 = c_n$  will guarantee a unique solution. Reinsch demonstrated the former, in which case these constraints imply  $\mathbf{c} = \mathbf{T}\mathbf{a}$ , where  $\mathbf{T}$  is an  $(n-2) \times n$  matrix,  $\mathbf{c} = (c_2, c_3, \dots, c_{n-1})'$  and  $\mathbf{a} = (a_1, a_2, \dots, a_n)'$ . Also note that  $f''(x) = 2c_i + \frac{2(c_{i+1}-c_i)}{h_i}(x-x_i)$ , so  $\int_{x_1}^{x_n} f''(x)^2 dx = \sum_{i=1}^{n-1} \int_{x_i}^{x_{i+1}} f''(x)^2 dx = \sum_{i=1}^{n-1} \frac{4}{3}h_i(c_i^2 + c_i c_{i+1} + c_{i+1}^2)$  is a quadratic form of  $\mathbf{c}$ . All together, the penalty can be written in a quadratic form  $\mathbf{a}'\mathbf{K}\mathbf{a}$ , where  $\mathbf{K}$  is an  $n \times n$  positive-semidefinite matrix with rank  $(n-2)$ , called a *penalty matrix*. Particularly,  $\mathbf{K}$  is a banded matrix with five diagonals, and its elements only depend on  $h_i$ . Let  $\mathbf{y} = (y_1, y_2, \dots, y_n)'$  and  $\mathbf{W} = (w_1, w_2, \dots, w_n)'$ , then it is easy to see that the penalized least squares problem becomes

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} \{(\mathbf{y} - \mathbf{a})'\mathbf{W}(\mathbf{y} - \mathbf{a}) + \lambda \mathbf{a}'\mathbf{K}\mathbf{a}\}.$$

For any known  $\lambda$ , the solution is  $\hat{\mathbf{a}} = (\mathbf{W} + \lambda\mathbf{K})^{-1}\mathbf{W}\mathbf{y}$ . Other parameters  $\hat{\mathbf{b}}$ ,  $\hat{\mathbf{c}}$  and  $\hat{\mathbf{d}}$  can be then determined from their relationship with  $\hat{\mathbf{a}}$ . The resulting spline  $\hat{f}(x)$  with estimated parameters is called a *smoothing spline*.

While  $f(x)$  is constructed as piecewise polynomials, it can also be written as a linear combination of  $n$  basis. To see this, consider an out-of-sample prediction for  $f(\tilde{x}_j)$ , where  $\tilde{x}_j$  is not any of  $x_i$  (otherwise the predicted value is just  $a_i$ ), but still inside  $[x_1, x_n]^1$ . Define  $\mathbf{D}_j$  to be an  $n \times n$  diagonal matrix with  $\mathbf{D}_j(i, i) = \tilde{x}_j - x_i$ , and  $\boldsymbol{\delta}_j$  to be a length- $n$  vector with  $\boldsymbol{\delta}_j(s_j) = 1$  and all other entries being zeros. The position index  $s_j$  satisfies  $x_{s_j} \leq \tilde{x}_j < x_{s_j+1}$ . Then we can write  $f(\tilde{x}_j) = \boldsymbol{\delta}_j' \mathbf{a} + \boldsymbol{\delta}_j' \mathbf{D}_j \mathbf{b} + \boldsymbol{\delta}_j' \mathbf{D}_j^2 \mathbf{c} + \boldsymbol{\delta}_j' \mathbf{D}_j^3 \mathbf{d}$ . The linear constraints between parameters imply that there are linear transformations from  $\mathbf{b}$ ,  $\mathbf{c}$ ,  $\mathbf{d}$  to  $\mathbf{a}$ , thus the result can simplify to the following form that only depends on  $\mathbf{a}$ :

$$f(\tilde{x}_j) = (\mathcal{B}_1(\tilde{x}_j) \quad \mathcal{B}_2(\tilde{x}_j) \quad \dots \quad \mathcal{B}_n(\tilde{x}_j)) \mathbf{a} = \sum_{i=1}^n \mathcal{B}_i(\tilde{x}_j) a_i.$$

To interpret this, replace  $\tilde{x}_j$  with  $x$  so that  $f(x) = \sum_{i=1}^n \mathcal{B}_i(x) a_i$ . Now it is clear that  $\mathcal{B}_i(x)$  is the  $i^{\text{th}}$  *basis function* and  $\mathcal{B}_i(\tilde{x}_j)$  is its value at  $x = \tilde{x}_j$ . In realistic computation with finite number of data point  $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_m$ , such basis representation ends up with a matrix equation  $\mathbf{f} = \mathbf{B}\mathbf{a}$ , where  $\mathbf{f} = (f(\tilde{x}_1), f(\tilde{x}_2), \dots, f(\tilde{x}_m))'$  and  $\mathbf{B}$  is an  $m \times n$  matrix, with  $\mathbf{B}(j, i) = \mathcal{B}_i(\tilde{x}_j)$ . Depending on the use of such matrix, it may be termed differently. For example, in the above out-of-sample prediction problem, it is known as a *predictor matrix*, because it transforms coefficients  $\mathbf{a}$  to predicted values. If all the basis functions are evaluated at sampling points  $x_1, x_2, \dots, x_n$  instead, the resulting matrix is known as a *design matrix*. Obviously, for this smoothing spline, the design matrix is an identity matrix. It also means that the  $i^{\text{th}}$  basis function  $\mathcal{B}_i(x)$  has value 1 at the  $i^{\text{th}}$  knot  $x_i$ , but value 0 at all other knots. Such basis is known as *cardinal spline basis*.

Thanks to the sparsity of  $\mathbf{K}$ , computation of a smoothing spline is efficient with  $O(n)$  complexity. However, at least since Wahba (1980) and Parker and Rice (1985), it has been recognised that using an  $n$  dimensional basis representation is computationally wasteful for negligible statistical gain. Gu and Kim (2002), Hall and Opsomer (2005), Li and Ruppert (2008), Kauermann et al. (2009), Claeskens et al. (2009) and Wang et al. (2011) set up a spline with a smaller number of knots, say  $k$ , so that it is spanned by  $k$  basis:  $f(x) = \sum_{i=1}^k \mathcal{B}_i(x) a_i$ . Such spline is known as a *low-rank penalized regression spline*. They show that  $k$  only grows rather slowly with sample size  $n$  to achieve asymptotically equivalent results to using a full smoothing spline (e.g.  $k = O(n^{1/5})$  for a cubic spline under REML smoothness estimation).

<sup>1</sup>In fact, this is just a cubic spline interpolation problem over  $(x_1, a_1), (x_2, a_2), \dots, (x_n, a_n)$ .



To use a penalized regression spline, one has to place knots, and their locations may not be at known sampling points. So in the basis representation of such a spline, basis coefficients only mean spline's values at knots not at sampling points. Another implication is that the design matrix is no longer an  $n \times n$  identity matrix, but an  $n \times k$  matrix. The penalty matrix is no longer an  $n \times n$  banded matrix, but a  $k \times k$  dense matrix. Let us respectively denote this new design matrix and new penalty matrix by  $\mathbf{X}$  and  $\mathbf{S}$ , the associated penalized least squares problem takes the following form:

$$\hat{\beta}_\lambda = \arg \min_{\beta} \{ \|\mathbf{W}^{\frac{1}{2}}(\mathbf{y} - \mathbf{X}\beta)\|^2 + \lambda \beta' \mathbf{S} \beta \}.$$

I specially index the estimated coefficients by  $\lambda$  to remind you that it is conditional on given  $\lambda$ . Later in §1.1.6 and §4.3 I will introduce the estimation of  $\lambda$ .

A low-rank cubic regression spline is not a sparse smoother, however. If we represent a spline as a B-spline (de Boor, 1978), whose basis function has *compact support*, we can obtain a sparse design matrix. In penalized regression, a *difference penalty matrix* is often used to preserve such sparsity in estimation. Furthermore, when knots are evenly spaced, such penalty has a simple interpretation as well. This type penalized regression spline is known as *P-splines* (Eilers and Marx, 1996; Marx and Eilers, 1998; Ruppert et al., 2003). Later, Wood (2011, §5.1) developed its adaptive variant, where the smoothing parameter is also represented by a B-spline, so that the strength of the penalization can vary between sampling locations. Recently, Wood (2017) also implemented the conventional derivative-based penalty (like  $\int_{x_1}^{x_n} f''(x)^2 dx$ ) with B-splines.

### 1.1.2 Low-rank penalized thin-plate regression spline

A smoothing problem is also defined in  $d > 1$  dimensions:

$$y_i = f(x_1, x_2, \dots, x_d) + \epsilon_i, \quad \epsilon_i \sim N(0, \phi/w_i),$$

where the smooth function  $f$  can be represented by a thin-plate spline or other Duchon spline (Duchon, 1977). For example, consider a bivariate example. Given  $n$  observations  $(x_{11}, x_{12}; y_1), (x_{21}, x_{22}; y_2), \dots, (x_{n1}, x_{n2}; y_n)$ , a thin-plate spline is the solution to the following penalized least squares problem:

$$f(x_1, x_2) = \arg \min_{f(x_1, x_2)} \left\{ \sum_{i=1}^n w_i (y_i - f(x_{i1}, x_{i2}))^2 + \lambda \iint \left( \frac{\partial^2 f}{\partial x_1^2} \right)^2 + 2 \left( \frac{\partial^2 f}{\partial x_1 \partial x_2} \right)^2 + \left( \frac{\partial^2 f}{\partial x_2^2} \right)^2 dx_1 dx_2 \right\},$$

and it has a direct basis representation:

$$f(x_1, x_2) = \sum_{i=1}^n \delta_i \mathcal{B}_i(x_1, x_2) + \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2,$$

where  $\mathcal{B}_i(x_1, x_2) = r_i^2 \log(r_i)$  with Euclidean distance  $r_i = \sqrt{(x_1 - x_{i1})^2 + (x_2 - x_{i2})^2}$  is a radial basis centred at a knot  $(x_{i1}, x_{i2})$ . Define an  $n \times n$  matrix  $\mathbf{E}$  with  $\mathbf{E}(j, i) = \mathcal{B}_i(x_{j1}, x_{j2})$  (note that this matrix is symmetric) and an  $n \times 3$  matrix  $\mathbf{T}$  with  $\mathbf{T}(i, ) = (1, x_{i1}, x_{i2})$ , then the penalized least squares can be shown to be (Wahba, 1990; Green and Silverman, 1994)

$$(\hat{\delta}, \hat{\alpha}) = \arg \min_{(\delta, \alpha)} \{ \|\mathbf{W}^{\frac{1}{2}}(\mathbf{y} - \mathbf{E}\delta - \mathbf{T}\alpha)\|^2 + \lambda \delta' \mathbf{E} \delta \}, \quad \text{subject to } \mathbf{T}' \delta = \mathbf{0},$$

where  $\mathbf{y} = (y_1, y_2, \dots, y_n)'$ ,  $\delta = (\delta_1, \delta_2, \dots, \delta_n)'$ ,  $\alpha = (\alpha_0, \alpha_1, \alpha_2)'$ .

Since  $\mathbf{E}$  is a dense matrix, computation for a thin-plate spline involves  $O(n^3)$  FLOP which is far too expensive even for moderately large  $n$ . Wood (2003) considered an eigen decomposition (see §4.1.5 if you need a revision) of  $\mathbf{E} = \mathbf{U} \mathbf{D} \mathbf{U}'$ , then by retaining the first  $k$  eigenvalues  $\mathbf{D}_k = \mathbf{D}(1:k, 1:k)$  and eigenvectors  $\mathbf{U}_k = \mathbf{U}(1:k, )$ , a low-rank approximation  $\mathbf{E}_k \approx \mathbf{U}_k \mathbf{D}_k \mathbf{U}_k'$  is obtained. Using  $\mathbf{E}_k$ , the penalized least squares problem becomes

$$(\hat{\delta}_k, \hat{\alpha}) = \arg \min_{(\delta_k, \alpha)} \{ \|\mathbf{W}^{\frac{1}{2}}(\mathbf{y} - \mathbf{U}_k \mathbf{D}_k \delta_k - \mathbf{T}\alpha)\|^2 + \lambda \delta_k' \mathbf{D}_k \delta_k \}, \quad \text{subject to } (\mathbf{U}_k' \mathbf{T})' \delta_k = \mathbf{0},$$

where  $\delta_k = U_k' \delta$  is a new length- $k$  parameter vector associated with a new  $n \times k$  design matrix  $U_k D_k$  and a new  $k \times k$  diagonal penalty matrix  $D_k$ . Solving this new problem only requires  $O(nk^2)$  FLOP. At first glance, one may wonder why this is beneficial. Given that the initial eigen decomposition has  $O(n^3)$  (or  $O(n^2k)$  if using *Lanczos algorithm*) complexity, the overall complexity for solving this penalized least squares problem is not any lower. However, in practice, the smoothing parameter  $\lambda$  needs be estimated, too. Many trial  $\lambda$  would be attempted, and a penalized least squares problem needs be solved for each trial. Note that the initial eigen decomposition only needs be done once, so all subsequent penalized least squares estimations can be  $O(n^2/k^2)$  times faster<sup>2</sup>.

For later out-of-sample prediction after estimation, a predictor matrix can be constructed; but to use it, the original parametrization is required. This follows from  $\delta = U_k \delta_k$ , or  $\begin{pmatrix} \delta \\ \alpha \end{pmatrix} = \begin{pmatrix} U_k \\ I \end{pmatrix} \begin{pmatrix} \delta_k \\ \alpha \end{pmatrix}$ .

### 1.1.3 Low-rank penalized tensor product regression spline

A thin-plate spline is isotropic due to its representation via radial basis. While this is useful for some settings, like spatial smoothing, it is not in others. Its basic problem is that it is not invariant to linear scaling of variables. For example, if we set  $\tilde{x}_2 = 2x_2$  in for the example bivariate thin-plate spline in the previous section, the estimated thin-plate spline will be different.

One way to construct a multivariate spline that is invariant to linear scaling, is via a *tensor product spline*. Such spline can be motivated in several different ways, but the most intuitive illustration may be via a *varying coefficient model*. Consider a bivariate smoothing problem again. We start from an obviously wrong smoothing problem  $y_i = g(x_{i1}) + \epsilon$  where we can obtain a cubic spline  $\hat{g}(x_1) = \sum_s \mathcal{B}_s(x_1) \alpha_s$ . Now to model the effect along the second dimension, we allow the coefficient  $\alpha_s$  to vary smoothly with  $x_2$ , giving another cubic spline  $\alpha_s(t) = \sum_t \mathcal{C}_t(x_2) \beta_{st}$ . Note that the new coefficient  $\beta_{st}$  is indexed by two subscripts, because the extension is made for each  $s$ . Now if we plug in  $\alpha_s(t)$  into  $\hat{g}$ , we get a new bivariate spline  $f(x_1, x_2) = \sum_s \sum_t \mathcal{B}_s(x_1) \mathcal{C}_t(x_2) \beta_{st}$ . It is clear that the basis for such spline is generated by all pairwise product of two sets of *marginal basis*, or, their *tensor product*, and the construction generalizes to any higher dimension.

When it comes to matrix or vector representation, computation of tensor product basis becomes computation of *Kronecker product*, denoted by  $\otimes$ . Specifically, the Kronecker product between vectors  $\mathbf{a} = (a_1 \ a_2 \ \cdots \ a_s)$  and  $\mathbf{b} = (b_1 \ b_2 \ \cdots \ b_t)$  is

$$\mathbf{a} \otimes \mathbf{b} = (a_1 \mathbf{b} \ a_2 \mathbf{b} \ \cdots \ a_s \mathbf{b}),$$

i.e,  $\mathbf{b}$  is distributed to each element of  $\mathbf{a}$  for multiplication. The resulting vector has length  $st$ . Now, suppose that given  $n$  observations, the marginal cubic spline basis for the first and second dimensions respectively end up with a *marginal design matrix*  $\mathbf{A}$  and  $\mathbf{B}$  of dimension  $n \times t$  and  $n \times s$ , then if  $\mathbf{C}$  is the design matrix for the tensor product basis, we have  $\mathbf{C}(i, ) = \mathbf{A}(i, ) \otimes \mathbf{B}(i, )$ . In other words, we are computing a row-wise Kronecker product between  $\mathbf{A}$  and  $\mathbf{B}$ , denoted by  $\mathbf{A} \tilde{\otimes} \mathbf{B}$  for convenience. The resulting *tensor product design matrix* has dimension  $n \times st$ .

In practice, a tensor product design matrix is never computed by row, as an alternative definition exists. To see this, consider an example with two  $2 \times 2$  matrices:

$$\begin{aligned} \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \tilde{\otimes} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} &= \begin{pmatrix} a_{11}b_{11} & a_{11}b_{12} & a_{12}b_{11} & a_{12}b_{12} \\ a_{21}b_{11} & a_{21}b_{12} & a_{22}b_{11} & a_{22}b_{12} \\ a_{11}b_{21} & a_{11}b_{22} & a_{12}b_{21} & a_{12}b_{22} \\ a_{21}b_{21} & a_{21}b_{22} & a_{22}b_{21} & a_{22}b_{22} \end{pmatrix} \\ &= \left( \begin{bmatrix} a_{11} & & & \\ & a_{21} & & \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \begin{bmatrix} a_{12} & & & \\ & a_{22} & & \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \right). \end{aligned}$$

In general, there is

$$\mathbf{A} \tilde{\otimes} \mathbf{B} = (\text{diag}(\mathbf{A}(, 1))\mathbf{B} \ \text{diag}(\mathbf{A}(, 2))\mathbf{B} \ \cdots \ \text{diag}(\mathbf{A}(, s))\mathbf{B}), \quad (1.1)$$

<sup>2</sup>If we further apply the model matrix reduction to be introduced in §4.2, we can further reduce the computational complexity of the penalized least squares problem to  $O(k^3)$ , hence the estimation process will be  $O(n^3/k^3)$  times faster.

---

```

 $p = \prod_{i=1}^t p_{B_i}$ 
allocate storage for  $B$ 
 $B = \text{zeros}(n, p)$ 
copy the last margin into the trailing block of  $B$ 
 $q = p_{B_t}$ 
 $B(:, (p - q + 1) : p) = B_t$ 
use backward recursion to fill in all entries of  $B$ 
for  $i = (t - 1) : 1$ 
     $r = q * p_{B_i}$ 
     $B(:, (p - r + 1) : p) = B_i \tilde{\otimes} B(:, (p - q + 1) : p)$ 
     $q = r$ 

```

---

**Figure 1.2:** An efficient algorithm to compute  $B = B_1 \tilde{\otimes} B_2 \tilde{\otimes} \cdots \tilde{\otimes} B_t$ .

i.e.,  $B$  is distributed to each column of  $A$  for row-scaling.

Kronecker product follows law of associativity. For three vectors  $\mathbf{a}$ ,  $\mathbf{b}$  and  $\mathbf{c}$ , this implies  $\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} = (\mathbf{a} \otimes \mathbf{b}) \otimes \mathbf{c} = \mathbf{a} \otimes (\mathbf{b} \otimes \mathbf{c})$ . For matrices  $A$ ,  $B$  and  $C$ , there is  $A \tilde{\otimes} B \tilde{\otimes} C = (A \tilde{\otimes} B) \tilde{\otimes} C = A \tilde{\otimes} (B \tilde{\otimes} C)$ . This offers an efficient algorithm for computing  $B = B_1 \tilde{\otimes} B_2 \tilde{\otimes} \cdots \tilde{\otimes} B_t$ , as presented in Figure 1.2, where  $B_i$  has dimension  $n \times p_{B_i}$ .

For smoothing purpose, a tensor product spline needs be penalized. Eilers and Marx (2003) used a weighted *Kronecker sum* of *marginal penalty matrices*. For example, if there are three margins, each with a penalty matrix  $K_i$ , then the tensor product penalty matrix is  $K = \lambda_1(K_1 \otimes I_2 \otimes I_3) + \lambda_2(I_1 \otimes K_2 \otimes I_3) + \lambda_3(I_1 \otimes I_2 \otimes K_3)$ , where the identity matrix  $I_i$  has the same dimension with  $K_i$ . Later Wood (2006) suggested that when possible (for example, if the margins are cubic regression splines not thin-plate regression splines) a reparametrization to each margins will give such penalty matrix a better interpretation. Note that the product here is the usual Kronecker product between matrices. The following is an explanatory example for its computation:

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \otimes B = \begin{pmatrix} a_{11}B & a_{12}B \\ a_{21}B & a_{22}B \end{pmatrix}$$

When marginal basis are low-rank spline basis, the resulting tensor product spline basis is also low-rank. Such spline basis is useful for modelling “interaction” between smooth functions. Examples include Belitz and Lang (2008), Augustin et al. (2009), Lee. and Durbán (2011) and Wood et al. (2013).

#### 1.1.4 Canonical form of a low-rank penalized regression spline

Splines introduced so far are those to be applied in this thesis, but there are also other variants of low-rank penalized regression splines. These include soap film smooths (Wood et al., 2008), splines on the sphere (Wahba, 1981) and Gaussian process smoothers (Kammann and Wand, 2003; Handcock et al., 1994), and they are useful in certain spatial applications. A common characteristic of all these low-rank penalized regression splines is that their associated smoothing problem has a canonical form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}'\mathbf{S}_\lambda\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{W}^{-1}\phi), \quad (1.2)$$

where  $\mathbf{X}$  is a design matrix,  $\boldsymbol{\beta}$  is a vector of coefficients,  $\mathbf{S}_\lambda$  is a penalty matrix (absorbing smoothing parameters),  $\boldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)'$  and  $\mathbf{W}$  is the inverse correlation matrix with  $\mathbf{W}(i, i) = w_i$  and  $\mathbf{W}(i, j) = 0, \forall i \neq j$ . The number of rows of  $\mathbf{X}$  is the number of data. The number of columns of  $\mathbf{X}$  gives the *rank* of the spline, and  $\mathbf{X}(:, i)$  gives the  $i^{\text{th}}$  spline basis evaluated at sampling points.

Clearly, a cubic regression spline is readily formulated as such. For a tensor product spline, it is also straightforward.  $\mathbf{X}$  is the tensor product design matrix, and the penalty matrix can be written as

$\mathbf{S}_\lambda = \sum_i \lambda_i \mathbf{S}_i$ , if we for example, define  $\mathbf{S}_1 = \mathbf{K}_1 \otimes \mathbf{I}_2 \otimes \mathbf{I}_3$  for the three-margin case used in the previous section.

For a thin-plate regression spline, some additional paper work is needed. First, we need to remove the linear constraints. These constraints imply that  $\delta_k$  is in the null space of  $(\mathbf{U}_k' \mathbf{T})'$ . From §4.1.1, we know that following a thin QR factorization  $\mathbf{U}_k' \mathbf{T} = \mathbf{Q} \mathbf{R}$ ,  $\mathbf{Q}^\perp$  is an orthonormal basis for such null space. Thus via another reparametrization  $\delta_k = \mathbf{Q}^\perp \delta_k^*$ , the penalized least squares problem becomes unconstrained in the new parameter space:

$$(\hat{\delta}_k^*, \hat{\alpha}) = \arg \min_{(\delta_k^*, \alpha)} \{ \|\mathbf{W}^{\frac{1}{2}}(\mathbf{y} - \mathbf{U}_k \mathbf{D}_k \mathbf{Q}^\perp \delta_k^* - \mathbf{T} \alpha)\|^2 + \lambda \delta_k^{*'} \mathbf{Q}^{\perp'} \mathbf{D}_k \mathbf{Q}^\perp \delta_k^* \}.$$

Then the result follows if we define  $\mathbf{X} = (\mathbf{U}_k \mathbf{D}_k \mathbf{Q}^\perp \quad \mathbf{T})$ ,  $\mathbf{S}_\lambda = (\lambda \mathbf{Q}^{\perp'} \mathbf{D}_k \mathbf{Q}^\perp \quad \mathbf{0})$  and  $\beta = (\delta_k^*, \alpha)'$ . Matrix multiplication with  $\mathbf{Q}^\perp$  only needs three Householder transformations (see §4.1.1 if you need a revision) hence is very efficient.

### 1.1.5 Bayesian interpretation of smoothing

The penalization on  $\beta$  can be interpreted by an (improper) prior  $\beta \sim N(\mathbf{0}, \mathbf{S}_\lambda^{-1} \phi)$  (Kimeldorf and Wahba, 1970; Wahba, 1983; Silverman, 1985; Fahrmeir and Lang, 2001; Ruppert et al., 2003), where  $\mathbf{S}_\lambda^{-1}$  is a *Moore-Penrose pseudoinverse* (recalling that  $\mathbf{S}_\lambda$  is rank-deficient). Then with the basis representation of a penalized regression spline, the canonical smoothing problem (1.2) is just a penalized linear model (or Bayesian linear model, linear mixed model):

$$\mathbf{y} | \beta \sim N(\mathbf{X} \beta, \mathbf{W}^{-1} \phi), \quad \beta \sim N(\mathbf{0}, \mathbf{S}_\lambda^{-1} \phi). \quad (1.3)$$

Following the basic Bayesian estimation procedure, we write down the (unnormalized) posterior of  $\beta$ :

$$P(\beta | \mathbf{y}) \propto P(\mathbf{y} | \beta) P(\beta) \propto \exp \left\{ - \frac{\|\mathbf{W}^{\frac{1}{2}} \mathbf{y} - \mathbf{W}^{\frac{1}{2}} \mathbf{X} \beta\|^2 + \beta' \mathbf{S}_\lambda \beta}{2\phi} \right\}.$$

For a known  $\lambda$ , the posterior mode is the solution to a *penalized weighted least squares problem*

$$\hat{\beta}_\lambda = \arg \min_{\beta} \|\mathbf{W}^{\frac{1}{2}} \mathbf{y} - \mathbf{W}^{\frac{1}{2}} \mathbf{X} \beta\|^2 + \beta' \mathbf{S}_\lambda \beta. \quad (1.4)$$

To find its solution, define the *weighted residual sum of squares* and the *penalized weighted residual sum of squares* by

$$D(\beta) = \|\mathbf{W}^{\frac{1}{2}} \mathbf{y} - \mathbf{W}^{\frac{1}{2}} \mathbf{X} \beta\|^2, \quad D_p(\beta, \lambda) = D(\beta) + \beta' \mathbf{S}_\lambda \beta, \quad (1.5)$$

then solve  $\partial D_p / \partial \beta = \mathbf{0}$ . It is easy to derive that

$$\begin{aligned} \frac{\partial D}{\partial \beta} &= -2 \mathbf{X}' \mathbf{W} (\mathbf{y} - \mathbf{X} \beta), & \frac{\partial^2 D}{\partial \beta \partial \beta'} &= 2 \mathbf{X}' \mathbf{W} \mathbf{X}, \\ \frac{\partial D_p}{\partial \beta} &= -2 (\mathbf{X}' \mathbf{W} \mathbf{y} + \mathbf{H}_\lambda \beta), & \frac{\partial^2 D_p}{\partial \beta \partial \beta'} &= 2 \mathbf{H}_\lambda, \end{aligned} \quad (1.6)$$

where

$$\mathbf{H}_\lambda = \mathbf{X}' \mathbf{W} \mathbf{X} + \mathbf{S}_\lambda. \quad (1.7)$$

Thus the solution to (1.4) is

$$\hat{\beta}_\lambda = \mathbf{H}_\lambda^{-1} \mathbf{X}' \mathbf{W} \mathbf{y}. \quad (1.8)$$

In fact, the posterior is a Gaussian distribution (because the exponent is a quadratic form of  $\beta$ ), so this posterior mode is also the posterior mean. Furthermore, the posterior variance is  $\mathbf{H}_\lambda^{-1} \phi$ , hence altogether, there is

$$\beta | \mathbf{y} \sim N(\hat{\beta}_\lambda, \mathbf{H}_\lambda^{-1} \phi).$$

Let  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}$  be fitted values and  $\hat{\mathbf{z}} = \mathbf{Z}\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}$  be some out-of-sample predicted values with predictor matrix  $\mathbf{Z}$ , then it is straightforward to see that they respectively have variance  $\mathbf{X}\mathbf{H}_{\boldsymbol{\lambda}}^{-1}\mathbf{X}'\phi$  and  $\mathbf{Z}\mathbf{H}_{\boldsymbol{\lambda}}^{-1}\mathbf{Z}'\phi$ . A convenient estimate (although not unbiased) of  $\phi$  is  $\hat{\phi} = D_p(\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}, \boldsymbol{\lambda}) / (n - \text{tr}(\mathbf{A}_{\boldsymbol{\lambda}}))$ , where  $\mathbf{A}_{\boldsymbol{\lambda}} = \mathbf{X}\mathbf{H}_{\boldsymbol{\lambda}}^{-1}\mathbf{X}'\mathbf{W}$  is a *hat matrix* or an *influence matrix* (so that  $\hat{\mathbf{y}} = \mathbf{A}_{\boldsymbol{\lambda}}\mathbf{y}$ ), and  $\text{tr}(\mathbf{A}_{\boldsymbol{\lambda}})$  is one definition of *effective degrees of freedom* for the estimated spline.

The above idea of estimation belongs to *empirical Bayesian inference*, not *hierarchical Bayesian inference* (or full Bayesian inference), as no prior distributions are assumed for  $\boldsymbol{\lambda}$  and  $\phi$ . However, given the link between smoothing and mixed modelling, several work (Zuur et al., 2014; Crainiceanu et al., 2005; Wood, 2016) has implemented full Bayesian inference via BUGS (Spiegelhalter et al., 1996) or JAGS (Plummer, 2003).

### 1.1.6 Smoothness selection via REML

Generally the smoothing parameter  $\boldsymbol{\lambda}$  needs be estimated. A conventional criterion is to choose one that minimizes the *generalized cross validation (GCV)* score (Craven and Wahba, 1979; Wood, 2004, 2008)

$$\text{GCV}(\boldsymbol{\lambda}) = \frac{D_p(\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}, \boldsymbol{\lambda})}{(1 - \text{tr}(\mathbf{A}_{\boldsymbol{\lambda}})/n)^2}.$$

However, Reiss and Ogden (2009) showed that a *restricted likelihood* (Patterson and Thompson, 1971) approach offers practical reliability advantages, being less prone to multiple local optima and consequent under-smoothing.

For the smoothing problem (1.2), the log restricted likelihood is

$$l_r(\boldsymbol{\lambda}, \phi) = \log \int \mathbf{P}(\mathbf{y}|\boldsymbol{\beta})\mathbf{P}(\boldsymbol{\beta})d\boldsymbol{\beta}.$$

The best estimate of  $\boldsymbol{\lambda}$  and  $\phi$  maximizes  $l_r(\boldsymbol{\lambda}, \phi)$ , or minimizes a *REML score*  $\mathcal{V}_r(\boldsymbol{\lambda}, \phi) = -2l_r(\boldsymbol{\lambda}, \phi)$  (Wood, 2011; Wood et al., 2015). The integral in  $l_r(\boldsymbol{\lambda}, \phi)$  has a closed form (as the integrand is the unnormalized posterior of  $\boldsymbol{\beta}$  which is Gaussian). Let  $m$  be the null space dimension of  $\mathbf{S}_{\boldsymbol{\lambda}}$  (or the number of  $\mathbf{S}_{\boldsymbol{\lambda}}$ 's zero eigenvalues) and  $|\mathbf{S}_{\boldsymbol{\lambda}}|_+$  be a generalized determinant (or the product of  $\mathbf{S}_{\boldsymbol{\lambda}}$ 's non-zero eigenvalues), then there are

$$\begin{aligned} \mathbf{P}(\mathbf{y}|\boldsymbol{\beta}) &= (2\pi\phi)^{-\frac{n}{2}}|\mathbf{W}|^{\frac{1}{2}}\exp\left\{-\frac{D(\boldsymbol{\beta})}{2\phi}\right\}, & \mathbf{P}(\boldsymbol{\beta}) &= (2\pi\phi)^{-\frac{p-m}{2}}|\mathbf{S}_{\boldsymbol{\lambda}}|_+^{\frac{1}{2}}\exp\left\{-\frac{\boldsymbol{\beta}'\mathbf{S}_{\boldsymbol{\lambda}}\boldsymbol{\beta}}{2\phi}\right\}, \\ \mathbf{P}(\mathbf{y}|\boldsymbol{\beta})\mathbf{P}(\boldsymbol{\beta}) &= (2\pi\phi)^{-\frac{n+p-m}{2}}|\mathbf{W}|^{\frac{1}{2}}|\mathbf{S}_{\boldsymbol{\lambda}}|_+^{\frac{1}{2}}\exp\left\{-\frac{D_p(\boldsymbol{\beta}, \boldsymbol{\lambda})}{2\phi}\right\}. \end{aligned}$$

Following an exact Taylor expansion of  $D_p(\boldsymbol{\beta}, \boldsymbol{\lambda})$  at its minimum point  $\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}$  (1<sup>st</sup> derivative at  $\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}$  is zero; see (1.6) for 2<sup>nd</sup> derivative):

$$D_p(\boldsymbol{\beta}, \boldsymbol{\lambda}) = D_p(\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}, \boldsymbol{\lambda}) + (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}})'\mathbf{H}_{\boldsymbol{\lambda}}(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}),$$

the REML score can be shown to be

$$\mathcal{V}_r(\boldsymbol{\lambda}, \phi) = (n - m) \log(2\pi\phi) + \frac{D_p(\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}, \boldsymbol{\lambda})}{\phi} + \log |\mathbf{H}_{\boldsymbol{\lambda}}| - \log |\mathbf{S}_{\boldsymbol{\lambda}}|_+ - \log |\mathbf{W}|. \quad (1.9)$$

The minimizer of  $\mathcal{V}_r$  needs be found numerically as there isn't an analytical solution. Later in §4.3 I will demonstrate the Newton-Raphson method for this minimization task in great details.

## 1.2 Additive models via penalized regression splines

An *additive model* is an extension of a smoothing problem. It has a structure like

$$y_i = f_1(x_i) + f_2(z_i) + f_3(u_{i1}, u_{i2}) + f_4(v_{i1}, v_{i2}, v_{i3}) + \cdots + \epsilon_i, \quad \epsilon_i \sim N(0, \phi/w_i),$$

where the observed data are modelled by a number of additive smooth functions and  $f_i$  is the  $i^{\text{th}}$  component. An additive model is a special case of a *generalized additive model (GAM)* (Hastie and Tibshirani, 1986, 1990) (described later in §4.4). In full generality, a combination of different smooth functions can be used for model representation. For example,  $f_1$  can be a kernel smoother,  $f_2$  can be a loess line,  $f_3$  can be a thin-plate spline, and  $f_4$  can be a tensor product spline with all margins being cubic splines. A convenient estimation method for such a model is via a *backing fitting* algorithm, implemented as an iterative, component-wise smoothing problem. Another appealing representation is to express all smooth functions by low-rank penalized regression splines (Wahba, 1990). Since all splines have a basis representation and a quadratic penalty, they can then be combined into a “larger” penalized regression problem. In this way, smoothness selection be carried out by minimizing GCV or REML. In addition, parametric effects and random effects can also be easily integrated into the model structure and be estimated.

Here is a simple, illustrative example for an additive model via penalized regression spline:

$$y_i = \alpha_0 + u_i\alpha_1 + v_i\alpha_2 + f_1(x_i) + f_2(z_{i1}, z_{i2}) + \epsilon_i, \quad \epsilon_i \sim N(0, \phi/w_i),$$

where  $\alpha_0 + u_i\alpha_1 + v_i\alpha_2$  is a parametric component via a linear regression on variables  $u$  and  $v$ , and  $f_1$  and  $f_2$  are some smooth functions represented with low-rank penalized regression splines. We assume that  $\mathbf{f}_1 = \mathbf{X}_1\boldsymbol{\beta}_1$  is a spline with a single penalty matrix  $\lambda_1\mathbf{S}_1$ , and  $\mathbf{f}_2 = \mathbf{X}_2\boldsymbol{\beta}_2$  is a spline with a total penalty matrix from  $q_2$  additive penalty matrices  $\mathbf{S}_2 = \sum_{i=1}^{q_2} \lambda_{2i}\mathbf{S}_{2i}$ . Following our previous introductions on penalized regression spline, we see that  $f_1$  may be a cubic spline, and  $f_2$  may be a tensor product spline. However, there could be other possibilities that I have mentioned. For example,  $f_1$  may also be a univariate thin-plate spline or a P-spline, and  $f_2$  may be an adaptive P-spline or any type of splines with additional penalization on its null space (Wood, 2011, §5.3). Given such setup, the additive model can then be written as

$$\mathbf{y} = \mathbf{X}_0\boldsymbol{\beta}_0 + \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \lambda_1\boldsymbol{\beta}_1'\mathbf{S}_1\boldsymbol{\beta}_1 + \boldsymbol{\beta}_2'\mathbf{S}_2\boldsymbol{\beta}_2 + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{W}^{-1}\phi), \quad (1.10)$$

where

$$\mathbf{X}_0 = \begin{pmatrix} 1 & u_1 & v_1 \\ 1 & u_2 & v_2 \\ \vdots & \vdots & \vdots \\ 1 & u_n & v_n \end{pmatrix}, \quad \boldsymbol{\beta}_0 = \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \end{pmatrix}$$

are the design matrix and the coefficient vector for the parametric component. By further defining

$$\mathbf{X} = (\mathbf{X}_0 \quad \mathbf{X}_1 \quad \mathbf{X}_2), \quad \boldsymbol{\beta} = \begin{pmatrix} \boldsymbol{\beta}_0 \\ \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{pmatrix}, \quad \mathbf{S}_\lambda = \begin{pmatrix} \mathbf{0} & & \\ & \lambda_1\mathbf{S}_1 & \\ & & \sum_{i=1}^{q_2} \lambda_{2i}\mathbf{S}_{2i} \end{pmatrix}, \quad (1.11)$$

the additive model (1.10) becomes a canonical smoothing model (1.2) with Bayesian interpretation (1.3) and penalized weighted least squares problem (1.4). Its GCV and REML scores also follow those for the smoothing model. For example, let  $m$  be the null space dimension of  $\mathbf{S}_\lambda$  (or the number of  $\mathbf{S}_\lambda$ 's zero eigenvalues), the REML score is just (1.9).

## Chapter 2

# Additive models with splines in practice: simulated examples

In this Chapter I will use two simulated examples to briefly demonstrate some practical aspects of statistical modelling via smoothing or additive models with splines. These include

- How to properly model data with a cyclic characteristic?
- How to place knots and choose appropriate number of knots for a spline?
- How to reasonably do exploratory data analysis on data with autocorrelation?
- How to estimate an additive model with AR(1) model error?

Simulated datasets are useful as they help assess how well a method work. In addition, I will cover some visualization methods for model inspection and checking. This Chapter is a “warm-up” for the the next Chapter on statistical modelling of real-world data.

### 2.1 Classic example with independent errors

Consider a signal function (i.e., the “truth”)

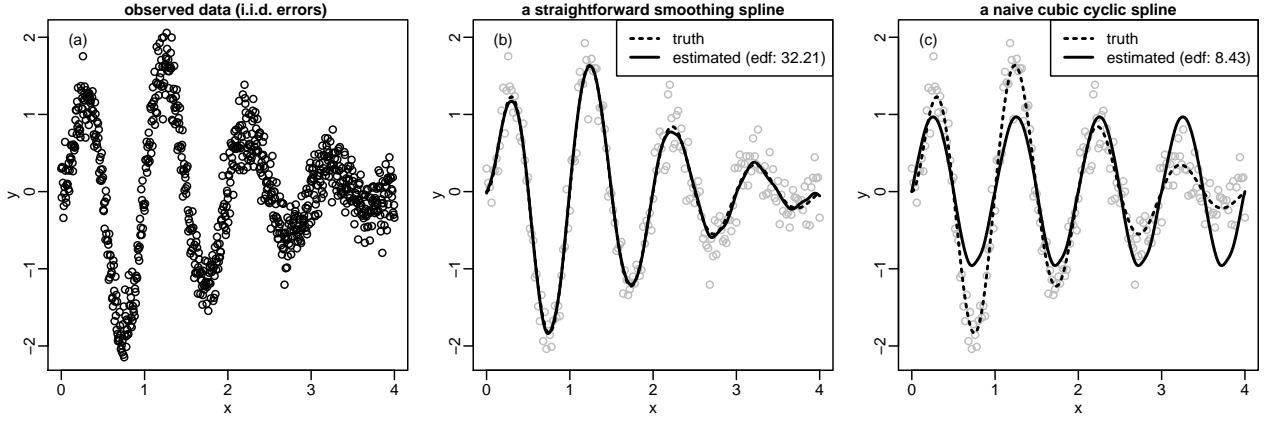
$$h(x) = 2\pi x \sin(2\pi x) e^{-\frac{2\pi x}{5}}, \quad x \in [0, 4].$$

This function factorizes into two parts:  $A(x) = 2\pi x e^{-\frac{2\pi x}{5}}$  and  $B(x) = \sin(2\pi x)$ . While  $h(x)$  is not a period function, it inherits the cyclic characteristic of  $B(x)$  and crosses zero at all whole numbers. 801 evenly spaced samples are taken on  $[0, 1]$ :  $x_i = 0.005(i - 1)$ ,  $i = 1, 2, \dots, 801$ , and assume that samples of the signal are observed with i.i.d. Gaussian noise

$$y_i = h(x_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \phi).$$

I will choose  $\phi = 0.09 \times \text{var}(h(x_i))$  so that the noise-to-signal ratio is 0.3. Panel (a) of Figure 2.1 sketches a random set of observed data (see Figure caption for the R code to reproduce them). Pretending that we don’t know the mathematical form of  $h(x)$ , we want to estimate it from the data using splines.

It is straightforward to fit a smoothing spline or a penalized regression spline on this dataset (see panel (b) of the Figure) and the fitted spline looks good, but it is undesired if we want to properly



**Figure 2.1:** Example data used for §2.1. The R code to reproduce them is `set.seed(0); x <- seq(0, 4, 0.005); A <- 2 * pi * x * exp(-2 * pi * x / 5); fx <- A * sin(2 * pi * x); y <- fx + rnorm(801, 0, 0.3 * sd(fx)); dat <- data.frame(x = x, A = A, fx = fx, y = y)`. Panel (a) is created by `with(dat, plot(x, y))`. The smoothing spline in panel (b) is obtained by `sm <- smooth.spline(dat$x, dat$y)`, where 140 knots (`sm$fit$nk - 2`) are placed by default. The resulting smoothing spline has degree of freedom 32.21 (obtained from `sm$df`). The simple cubic cyclic spline is estimated in `mgcv`, with `cycl <- gam(y ~ s(x, bs = "cc", k = 16), data = dat, knots = list(x = c(0, 1), method = "REML"))`. The effective degree of freedom is 8.43 (obtained from `sum(cycl$edf)`). `gam` is used here because using `bam` for such a small dataset is an overkill. Note that in panel (c) and (d), as well as all subsequent Figures associated with this example, only every 4 data are plotted (in gray dots), i.e., at  $i = 1, 5, 9, \dots, 801$ . This is for a clearer display of the true signal (dashed line) and the estimated spline (solid line).

model the cyclic effect (I will elaborate on this later in §2.1.2). By contrast, the cyclic effect can be easily captured by using a cubic cyclic spline (a variant of cubic regression spline introduced in §1.1.1, by modifying the natural boundary condition to the periodic condition where the function values equal at two ends), but the estimation result (see panel (c) of Figure 2.1) is not satisfying as it fails to model the variability in the data as  $x$  changes. Can we add more components onto the cyclic spline to get a more appropriate fit?

### 2.1.1 Tensor product spline for interaction

It is not obvious how this can be done. The classic “trend + seasonality” decomposition does not help here as the data clearly have no trend, leaving the “seasonality” as the only component.

To get some inspiration let us do some cheating. We know that  $h(x_i) = A(x_i)B(x_i)$  in the true signal. First of all, if we know every  $A_i = A(x_i)$ , estimation of  $B(x)$  should be very accurate from the following model

$$y_i = A_i f_0(x_i) + \epsilon_i,$$

where  $f_0(x)$  is a cubic cyclic spline to estimate  $B(x)$ . Note that the presence of  $A_i$  does not make model estimation any more difficult. By expressing the model in matrix form

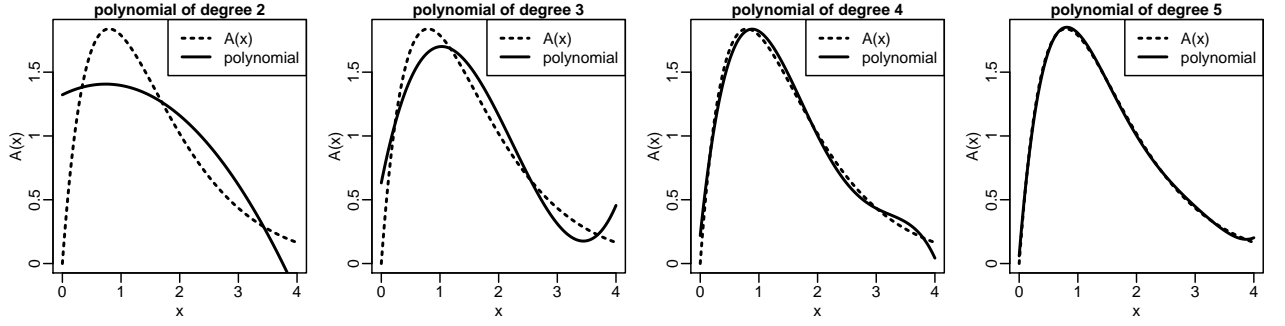
$$\mathbf{y} = \text{diag}(\mathbf{A})\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where  $\mathbf{A} = (A_1, A_2, \dots, A_n)$  collects all  $A_i$ ’s into a vector, we see that estimation of  $\boldsymbol{\beta}$  is just a regression problem with design matrix  $\text{diag}(\mathbf{A})\mathbf{X}$ .

Now what can we do if  $A_i$  is unknown? Can we approximate it with a polynomial of  $x_i$ , giving an approximation function

$$\left( \sum_{k=0}^p x_i^k \alpha_k \right) f_0(x_i) = f_0(x_i) + \sum_{k=1}^p x_i^k \alpha_k f_1(x_i) = f_0(x_i) + \sum_{k=1}^p x_i^k f_k(x_i),$$





**Figure 2.2:** Polynomial approximation to  $A(x)$ . The approximation is good enough when degree reaches 5. The regression polynomials are obtained via simple least squares fitting with function `lm`. For example, the cubic polynomial is obtained by `lm(A ~ poly(x, degree = 3), data = dat)`.

where the unknown polynomial coefficient  $\alpha_k$  is absorbed into the unknown cyclic spline  $f_0(x_i)$  for a new unknown cyclic spline  $f_k(x_i)$ ? As Figure 2.2 shows, this seems a reasonable idea. A polynomial of degree 5 is a very good approximation to  $A(x)$ .

We can also approximate  $A_i$  with a cubic regression spline  $\sum_s \mathcal{B}_s(x_i)\alpha_s$  rather than a high order polynomial (as said at the beginning of §1.1.1, this is desirable). If we express  $f_0(x_i) = \sum_t \mathcal{C}_t(x_i)\beta_t$ , we see that the approximation function

$$\left(\sum_s \mathcal{B}_s(x_i)\alpha_s\right)\left(\sum_t \mathcal{C}_t(x_i)\beta_t\right) = \sum_t \sum_s \mathcal{B}_s(x_i)\mathcal{C}_t(x_i)\alpha_s\beta_t := \sum_t \sum_s \mathcal{B}_s(x_i)\mathcal{C}_t(x_i)\gamma_{st}$$

is a tensor product spline, where one margin is a natural cubic spline of  $x$  and the other is a cubic cyclic spline of  $x$ .

You may have been confused why a tensor product spline, a multivariate spline, turns up in this univariate context. To clarify this, let us use variable  $x$  for the natural cubic spline and variable  $z = x - \lfloor x \rfloor$  for the cyclic cubic spline, where  $\lfloor x \rfloor$  is the biggest integer no larger than  $x$ . Note that  $x \in [0, 4]$ , but  $z \in [0, 1]$  (where a period of the cyclic spline is defined). Then we can express the true signal  $h(x)$  as a bivariate function  $h(x, z) = A(x)B(z)$ . From now on, I will use this bivariate representation for the true signal and approximation functions.

In summary, we have come up with the following candidate models to estimate  $h(x, z)$ :

**Model 2.1:** “true” model

$$y_i = A_i f_0(z_i) + \epsilon_i,$$

**Model 2.2:** polynomial model

$$y_i = f_0(z_i) + \sum_{k=1}^p x_i^k f_k(z_i) + \epsilon_i,$$

**Model 2.3:** tensor product spline

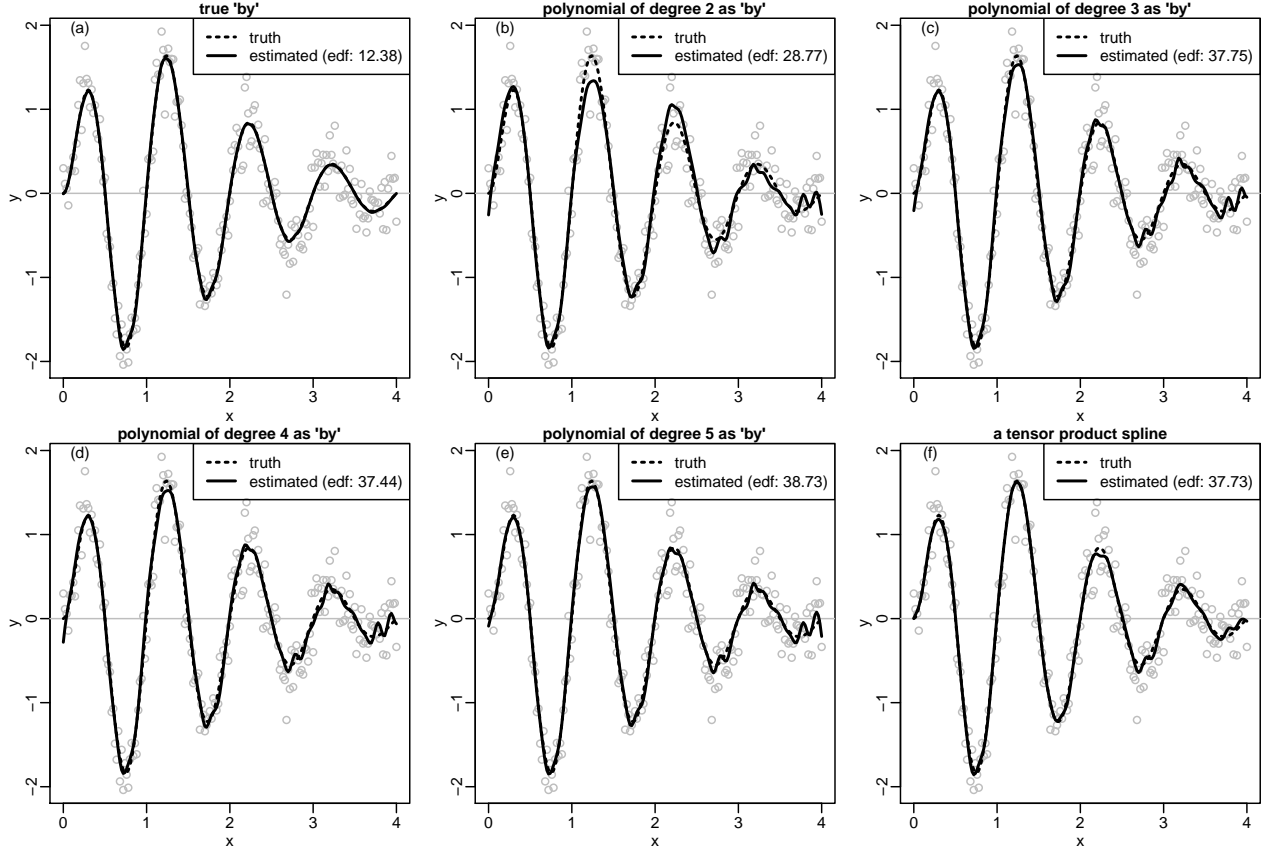
$$y_i = f_0(z_i) + f_1(z_i, x_i) + \epsilon_i.$$

I would like to emphasize that the tensor product spline  $f_1(z_i, x_i)$  should not be expressed in a separable form like  $g_1(z_i)g_2(x_i)$ , even if it has  $g_1$  and  $g_2$  as its margins. The separable form is a common misconception. In fact, a tensor product spline is much more flexible than the mere spline product. Suppose  $g_1$  and  $g_2$  respectively have  $k_1$  and  $k_2$  spline coefficients, the tensor product spline  $f_1$  will have  $k_1 k_2$  spline coefficients, while the spline product  $g_1 g_2$  only has  $k_1 + k_2$  spline coefficients.

All three models can be easily fitted with R package `mgcv`, using the REML estimation method introduced in §4.3. Both natural cubic splines and cubic cyclic splines are knots based splines so

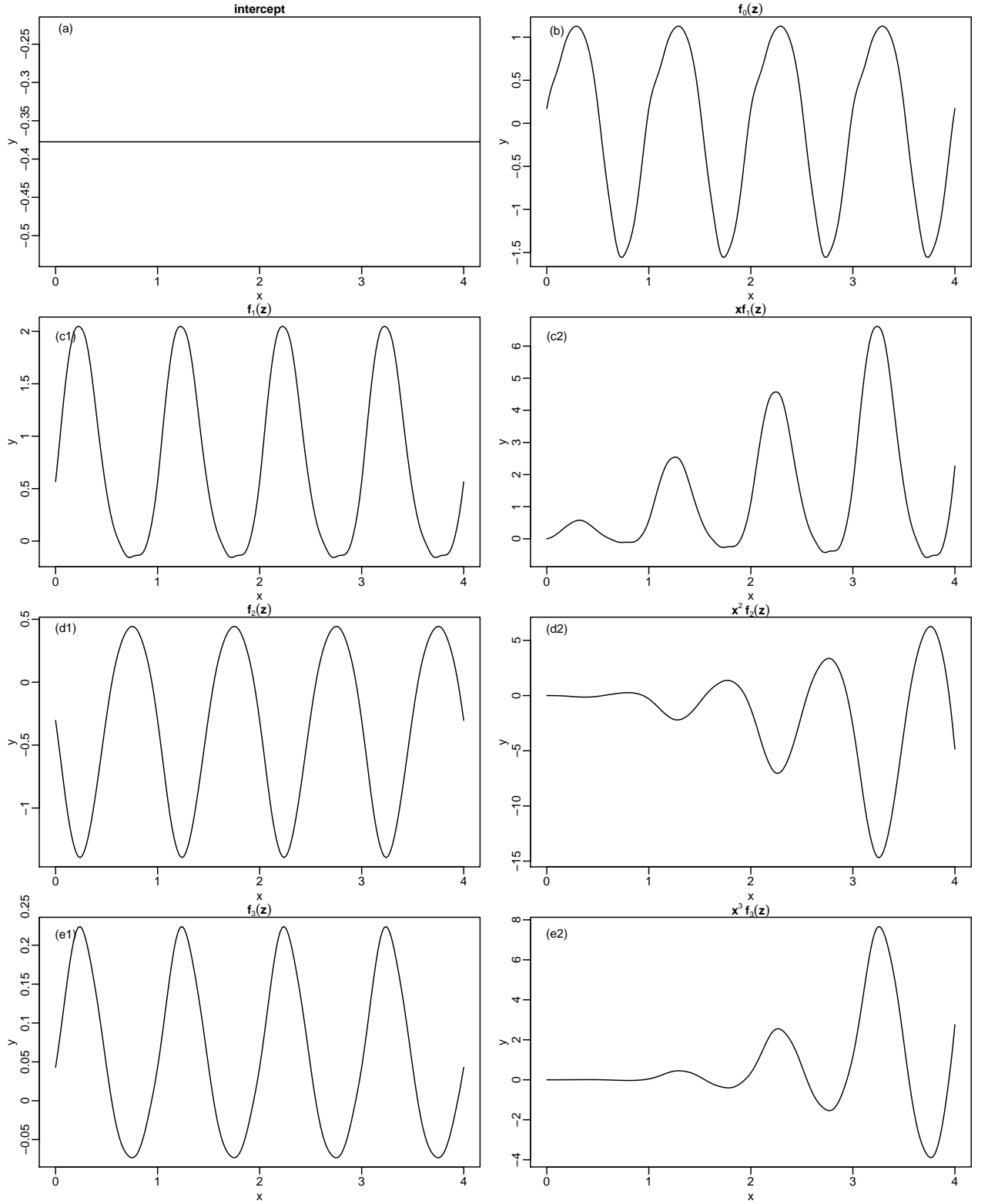
we have to supply knots so that spline basis can be constructed and evaluated to produce a design matrix. I will discuss knots related issues in the coming section, §2.1.3. For now, it is sufficient to know that setting 10 ~ 16 equally spaced knots for both  $x$  and  $z$  in their range are more than sufficient.

Figure 2.3 illustrates fitted splines for all these models. The best fit is surely model 2.1 (see panel (a)). No approximation for  $A(x)$  is done, so the cyclic spline can be accurately estimated. Nevertheless, this model is “cheating” as we should pretend that we don’t know  $A(x)$ . Panels (b) to (e) are model 2.2 with polynomials of degree 2, 3, 4 and 5. Although it appears from Figure 2.2 that polynomial approximation to  $A(x)$  is only good enough when degree reaches 5, estimation for  $h(x, z)$  is readily good enough when a cubic polynomial is used. Nevertheless, model 2.2 evidently overfits the data. Panel (f) sketches the fitted model 2.3. While there is also a little overfitting near their right boundary, it is less pronounced than that in the polynomial model.



**Figure 2.3:** Fitting models 2.1, 2.2 and 2.3 in `mgcv`. First create variable  $z$ : `dat$z <- dat$x - floor(dat$x)`. For the fitted spline in panel (a), use `gam(y ~ s(z, bs = "cc", k = 16, by = A), data = dat, knots = list(z = c(0, 1)), method = "REML")`. For the fitted spline in panel (c) for example, use `gam(y ~ s(z, bs = "cc", k = 16) + s(z, bs = "cc", k = 16, by = x) + s(z, bs = "cc", k = 16, by = I(x*x)), data = dat, knots = list(z = c(0, 1)), method = "REML")`. For the fitted tensor product spline in panel (f), use `gam(y ~ s(z, bs = "cc", k = 16) + ti(x, z, bs = c("cr", "cc"), k = c(12, 10)), data = dat, knots = list(z = c(0, 1)), method = "REML")`. Note the use of `ti` not `te` for properly creating interaction.

There is yet another reason why using polynomials should not be recommended, as is illustrated in Figure 2.4. The Figure plots every additive component in the model when the polynomial has degree 3. In other words, the fitted spline in panel (c) of Figure 2.3 are decomposed as (a) + (b) + (c2) + (d2) + (e2) of this Figure. It is concerning that terms in (c2), (d2) and (e2) cancel out each other. For example,  $xf_1(z)$  is positive in (roughly)  $[3, 3.7]$ , but  $x^2f_2(z)$  is negative in that range. This is an alternative illustration of Runge’s phenomenon (previously mentioned in §1.1.1), i.e., the oscillation of high order polynomials. In practice, such “cancellation” lacks numerical stability as many significant digits may be lost in the way, giving inaccurate computation result in finite-precision floating-point arithmetic.



**Figure 2.4:** Each additive component of model 2.2. Panel (a) is the intercept; Panel (b) is  $\hat{f}_0$  (the estimated function for  $f_0$ ). Panel (c1), (d1) and (e1) are  $\hat{f}_1$ ,  $\hat{f}_2$  and  $\hat{f}_3$ ; they are multiplied respectively by  $x$ ,  $x^2$  and  $x^3$  when they enter the model, which are sketched in panel (c2), (d2) and (e2). It is concerning that terms in (c2), (d2) and (e2) cancel out each other, as they have different sign in the same range. Such “cancellation” makes finite-precision numerical computation less accurate.

### 2.1.2 Why is an additive model better than a simple smoothing spline?

I have previously mentioned that using a smoothing spline to model  $h(x, z)$  is not a good idea; here is the reason. Suppose now our samples are not even spaced over  $[0, 4]$ . For simplicity let us assume that some of the original evenly spaced data are missing. Let us examine two cases:

1. Data near two boundaries are missing: observations indexed by  $i = 1, 2, \dots, 100$  and  $i = 702, 703, \dots, 801$  are missing;
2. Data in the middle are missing: observations indexed by  $i = 201, 201, \dots, 400$  are missing.

For each case we fit a simple smoothing spline and a tensor product model 2.3 and make prediction at those unsampled positions. The results are illustrated in Figure 2.5.

1. In the first case, prediction is an extrapolation problem. The natural boundary condition assumed by a smoothing spline leads to linear extrapolation (see panel (1a)). The tensor product spline successfully preserves the cyclic nature of the signal, although it overestimates the amplitude at the left boundary.
2. In the second case, prediction is an interpolation problem. The smoothing spline (see panel (2a)) underestimates the true function while the tensor product spline (see panel (2b)) yields good approximation to the truth.

### 2.1.3 Knots placement and choosing number of knots

In the previous section various splines have been fitted to the example data, but I haven't explained where, how and how many knots are placed for construction of these splines.

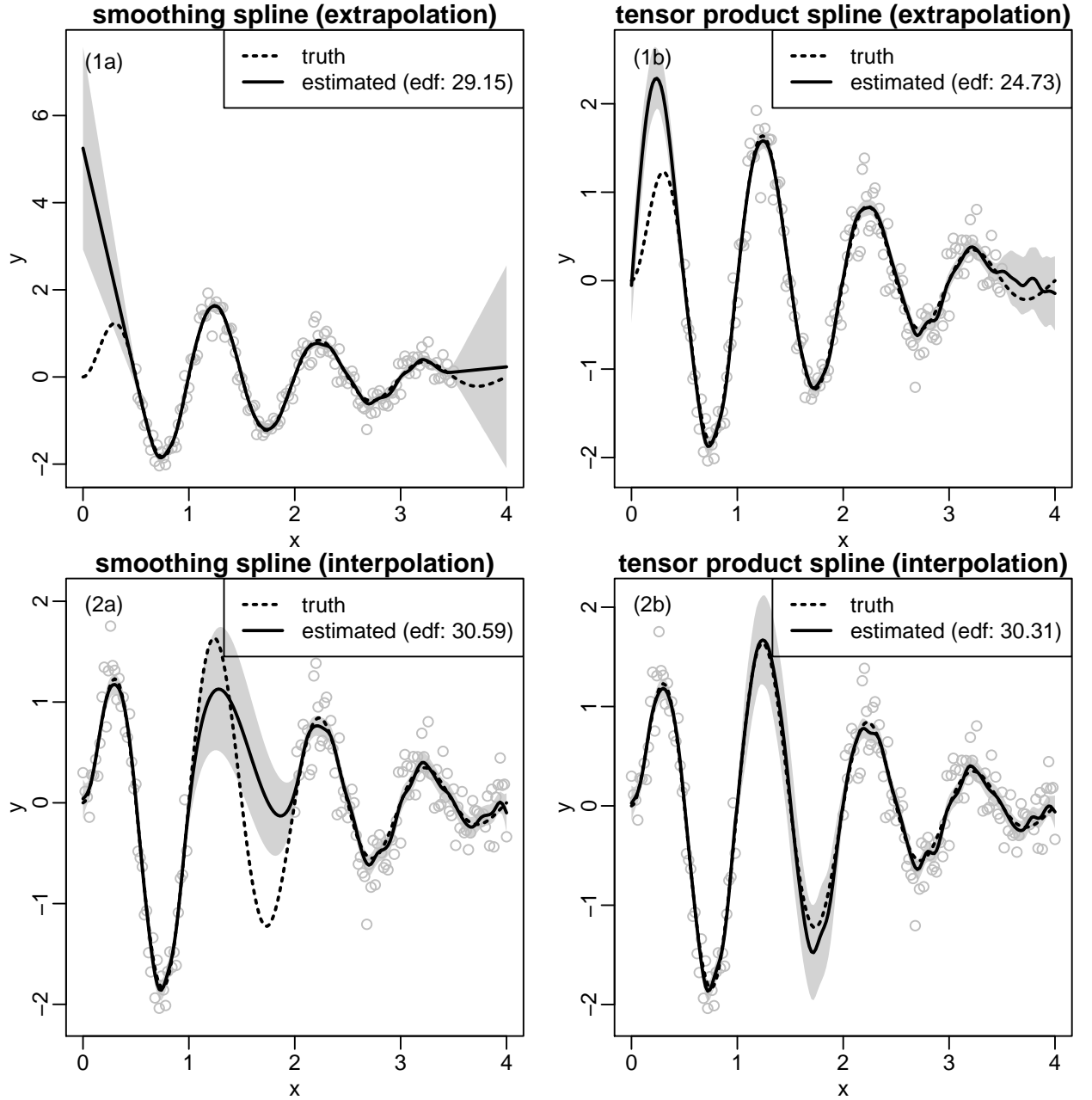
In general, suppose we want to place  $k$  knots  $g_1, g_2, \dots, g_k$  in order to estimate a regression spline  $y_i = f(x_i) + \epsilon_i$ , the basic requirements on knots placement are

1. At most one knot exist between any two adjacent unique data values  $x_i$  and  $x_{i+1}$ ;
2. At least one unique data value exists between any two adjacent knots  $g_j$  and  $g_{j+1}$

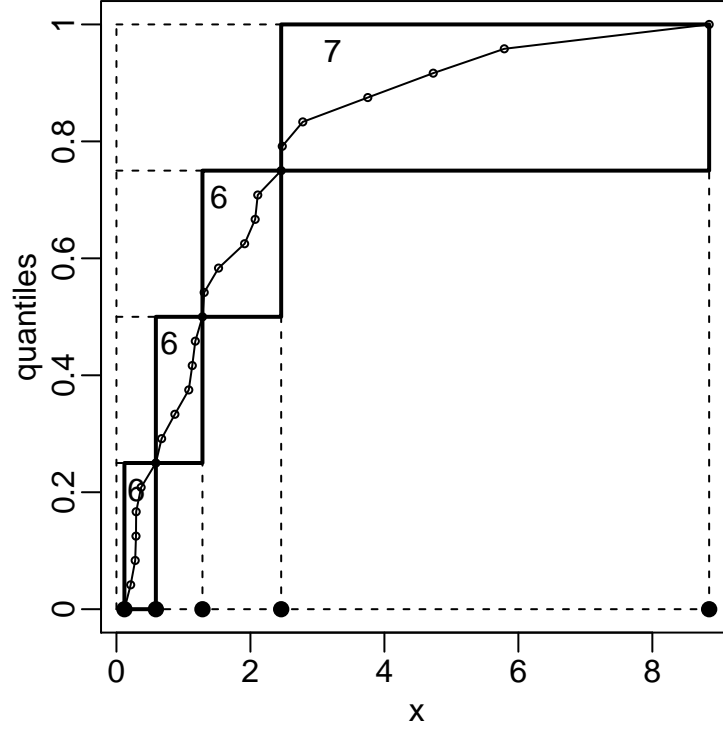
These two statements are in fact equivalent. Violation of either will lead to more than one knot between two two unique data values, which is problematic. For example, suppose there are two knots  $g_j$  and  $g_{j+1}$  between data values  $x_i$  and  $x_{i+1}$ , then the piecewise polynomial on  $[g_j, g_{j+1}]$  can not be uniquely determined as there are no data on this interval.

There is no unique way for knots placement, however, a canonical and automatic knots placement method is to place knots by quantiles of unique data values. This will meet both requirements, and achieve a stronger result on the second: there will be roughly equal number of unique data values on any  $[g_j, g_{j+1}]$ . See Figure 2.6 for an illustration.

If knots are automatically placed as above, knots placement is only up to choosing the number of knots. Choosing an appropriate number of knots is important for a regression spline, since the number of knots,  $k$ , decides the number of coefficients to parametrize a spline: the larger  $k$  is, the more flexible the spline is hence the closer it can approximate the data. When  $k$  is insufficient big, the spline can



**Figure 2.5:** Fitting a smoothing spline and a tensor product model 2.3 to non-evenly spaced samples on  $[0, 4]$ . In panel (1a) and (1b), data are missing near the two boundaries: `subset1 <- dat[101:701, ]`. The smoothing spline in panel (1a) is estimated by `gam(y ~ s(x, bs = "cr", k = 50), data = subset1, method = "REML")` and the tensor product spline in panel (1b) is estimated by `gam(y ~ s(xt, bs = "cc", k = 16) + ti(xt, x, bs = c("cc", "cr"), k = c(12, 10)), data = subset1, method = "REML")`. In panel (2a) and (2b), data are missing in the middle: `subset2 <- dat[c(1:200, 401:801), ]`. The smoothing spline in panel (2a) is estimated by `gam(y ~ s(xt, bs = "cc", k = 16) + ti(xt, x, bs = c("cc", "cr"), k = c(12, 10)), data = subset2, method = "REML")` and the the tensor product spline in panel (2b) is estimated by `gam(y ~ s(xt, bs = "cc", k = 16) + ti(xt, x, bs = c("cc", "cr"), k = c(12, 10)), data = subset2, method = "REML")`. 95% Bayesian confidence intervals (gray shaded) are produced by `plot.gam`. In both extrapolation and interpolation scenarios, the smoothing spline fails to well model the cyclic nature of the true signal.



**Figure 2.6:** Automatic knots placement by quantiles of unique data values. Suppose we have 25 unique  $x$ -values: `set.seed(0); n <- 25; x <- sort(rexp(n, 0.5))`, placing five knots by quantiles gives `k <- 5; g <- quantile(x, prob = seq(0, 1, length.out = k))`. The line plot is produced by lining up  $(x_i, \frac{i-1}{n-1})$ , i.e., `xp <- seq.int(0, n, length.out = n); plot(x, xp, type = "l"); points(x, xp, cex = 0.5)`. Consider slicing the vertical axis by `seq(0, 1, length.out = k)`, then the  $x$ -coordinates where these horizontal lines intersects the line plot are the quantiles (marked by bold solid dots on the  $x$ -axis). The slicing cuts the line into  $k - 1$  segments, and the graph above has boxed each piece by a rectangular. The number of data in each box is just the number of data between two adjacent knots. For this example there are 6, 6, 6, 7 data on the four segments (from left to right).

underfit the data; if  $k$  is too big, the spline can overfit the data. Ultimately when there are as many knots as data, we have an interpolation spline. The smoothing spline introduced in §1.1 does not end up with interpolation because of the penalization, which suppresses the effective degree of freedom of each spline coefficient to some values much smaller than one. A smoothing spline eliminates any knots related issues, but at the price of computational complexity. A low-rank penalized regression spline is less computationally expensive. Its basic idea is to practically choose a  $k$ , much smaller than the number of data, but big enough for the penalization to take effect.

To understand what this means, consider choosing  $k$  for our tensor product model 2.3. To facilitate the discussion here, I will rewrite it as

$$y_i = f_0(z_i; k_0) + f_1(z_i, x_i; k_{1,1}, k_{1,2}) + \epsilon_i,$$

where  $k_0$  is the number of knots for  $f_0$ ,  $k_{1,1}$  is the number of knots for  $z$  margin of  $f_1$  and  $k_{1,2}$  is the number of knots for  $x$  margin of  $f_1$ . Choosing those  $k$  values can be an iterative process. Here is an example demonstration.

1. Make an initial attempt with

$$y_i = f_0(z_i; 4) + f_1(z_i, x_i; 4, 3) + \epsilon_i,$$

then extract the partial residual w.r.t.  $f_0(z_i, 4)$  and fit

$$\text{partial residuals} = f_0(z_i, k_0) + \epsilon_i$$

for a set of increasingly big  $k_0$ . For each trial record the number of coefficients (*ncoef* for short) in  $f_0$  and the resulting effective degree of freedom (*edf* for short) of  $\hat{f}_0$ . The panel (a) of

Figure 2.7 sketches the *edf* against *ncoef*, with the corresponding  $k_0$  labelled above the points. When penalization has no effect, no coefficients are suppressed so we should have  $edf \approx ncoef$ . However, as  $k_0$  grows we observe that *edf* plateaus, because the penalty starts to play its role. From the graph it seems that  $k_0 = 16$  is an adequate choice.

2. Update the model with  $k_0 = 16$ , i.e., fit

$$y_i = f_0(z_i; 16) + f_1(z_i, x_i; 4, 3) + \epsilon_i,$$

then extract the partial residual w.r.t.  $f_1(z_i, x_i; 4, 3)$  and fit

$$\text{partial residuals} = f_1(z_i, x_i; k_{1,z}, 3) + \epsilon_i$$

for a set of increasing  $k_{1,1}$ . The panel (b) of Figure 2.7 sketches *edf* of  $\hat{f}_1$  against its *ncoef*. It appears that  $k_{1,1} = 12$  is sufficiently big.

3. Update the model with  $k_{1,1} = 12$  i.e., fit

$$y_i = f_0(z_i; 16) + f_1(z_i, x_i; 12, 3) + \epsilon_i,$$

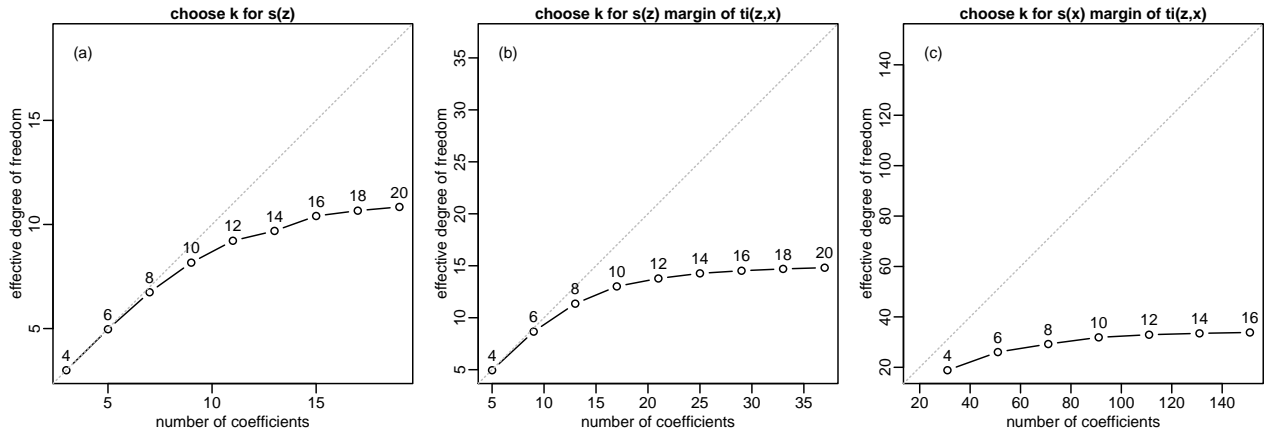
then extract the partial residual w.r.t.  $f_1(z_i, x_i; 12, 3)$  and fit

$$\text{partial residuals} = f_1(z_i, x_i; 12, k_{1,2}) + \epsilon_i$$

for a set of increasing  $k_{1,2}$ . The panel (c) of Figure 2.7 sketches *edf* of  $\hat{f}_1$  against its *ncoef*. It appears that  $k_{1,2} = 10$  is big enough.

So in the end we have a model

$$y_i = f_0(z_i; 16) + f_1(z_i, x_i; 12, 10) + \epsilon_i.$$

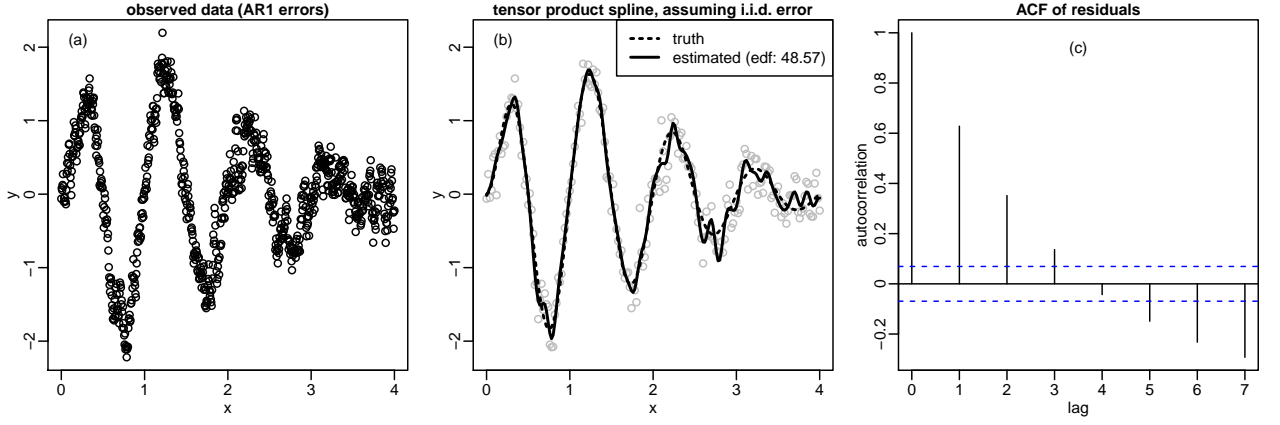


**Figure 2.7:** Choose the number of knots for model 2.3. These graphs illustrate the relationship between effective degree of freedom (*edf*, y-axis) and number of parameters (*ncoef*, x-axis), with the input  $k$  values labelled above the points. Due to penalization, The *edf* will not increase linearly with *ncoef*; it starts to plateau at some  $k$  value. This underlies the principle for choosing  $k$ . In practice an additive model has many spline terms so the selection process is a term-wise, iterative process. See main text for more explanations.

## 2.2 Time series example with AR(1) autocorrelation

This section adapts the example in §2.1 to a simple time series case study, replacing the i.i.d. error assumption with an AR(1) process with autocorrelation  $\omega = 0.75$ .

$$y_i = h(x_i) + e_i, \quad e_i = \omega \cdot e_{i-1} + \epsilon_i, \quad \epsilon_i \sim N(0, (1 - \omega^2)\phi).$$



**Figure 2.8:** Example data used for §2.2. The R code to reproduce them is `set.seed(0); omega <- 0.75; x <- seq(0, 4, 0.005); A <- 2 * pi * x * exp(-2 * pi * x / 5); fx <- A * sin(2 * pi * x); y <- fx + base::c(arima.sim(list(ar = omega), n, sd = sqrt(1 - omega * omega) * 0.3 * sd(fx))); dat <- data.frame(x = x, z = x - floor(x), fx = fx, y = y)`. Panel (a) is created by `with(dat, plot(x, y))`. The tensor product spline (see model 2.3) in panel (b) is obtained by `fit <- gam(y ~ s(z, bs = "cc", k = 16) + ti(x, z, bs = c("cr", "cc"), k = c(12, 10)), data = dat, knots = list(z = c(0, 1)), method = "REML")`, but it ends up overfitting the data. The ACF in panel (c) is generated by `acf(fit$resi, lag.max = 7)`.

$\text{var}(\epsilon_i) = (1 - \omega^2)\phi$  is specially chosen so that  $\text{var}(e_i)$  is still  $\phi$ , i.e., the noise-to-signal ratio is as same as the previous example. Panel (a) of Figure 2.8 sketches a random set of observed data, and see Figure caption for R code for reproducing it.

Autocorrelated error is known to cause overfitting in smoothing models, if errors are wrongly assumed to be i.i.d.. For example, see Krivobokova and Kauermann (2007); Opsomer et al. (2001); Smith et al. (1998); Wang (1998). Panel (b) of Figure 2.8 verifies this, where the tensor product spline model 2.3 assuming i.i.d. errors ends up overfitting the data. Panel (c) graphs the (sample autocorrelation function) ACF of model residuals. The residuals are still highly correlated even at presence of overfitting.

### 2.2.1 How to do exploratory analysis for data with autocorrelation

The existence of autocorrelation can complicate or even mislead model development. For example, let us consider choosing  $k$  for all splines in the tensor product model. In §2.1.3 it has been demonstrated that  $k_1 = 16$ ,  $k_{1,1} = 12$  and  $k_{1,2} = 10$  are readily sufficient to approximate the truth. However, when observations are autocorrelated, this “fact” can be obscured. Figure 2.9 illustrates what will happen if we replicate exactly the same selection process on  $k$  as described in §2.1.3. Panel (a) shows that  $k_1 = 16$  is far from being adequate; we probably need to choose  $k_1 = 60$ . Then in panel (b), the *edf* almost grows linearly and does not seem to plateau.

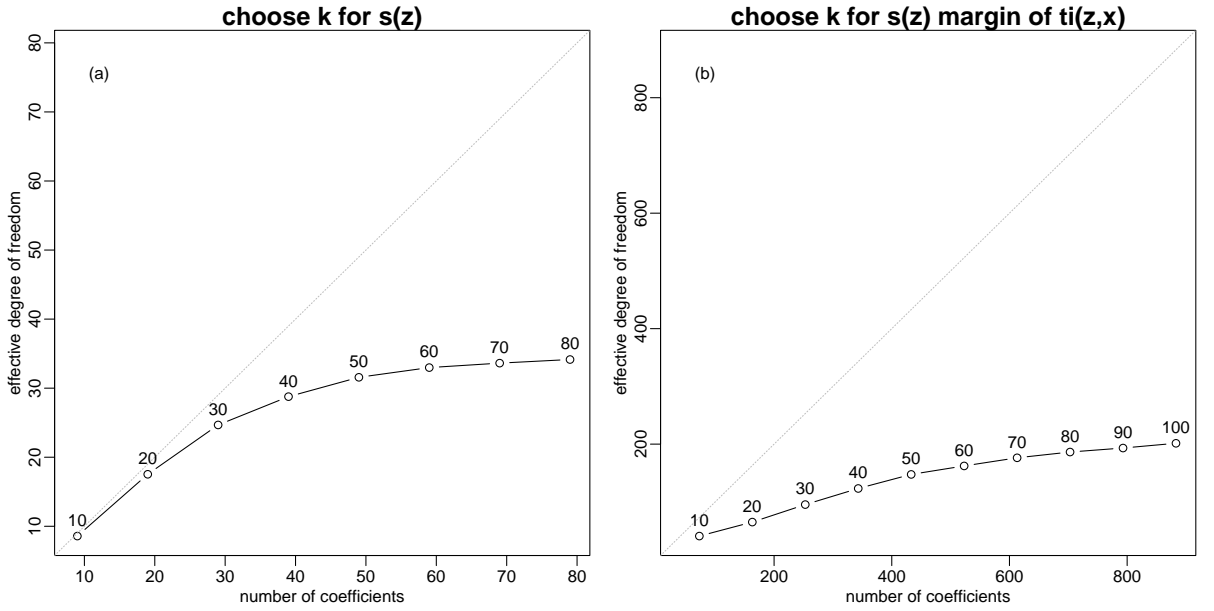
A useful trick for exploratory anaysis on data with autocorrelation is thinning. For example, let us examine every 4<sup>th</sup> data indexed  $i = 1, 5, 9, \dots, 801$ . Figure 2.10 shows that  $k$  can now be properly chosen. Figure 2.11 further shows that fitting the tensor product model 2.3 on the thinned data yields less misleading estimation.

### 2.2.2 AR(1) error with autocorrelation coefficient $\omega$

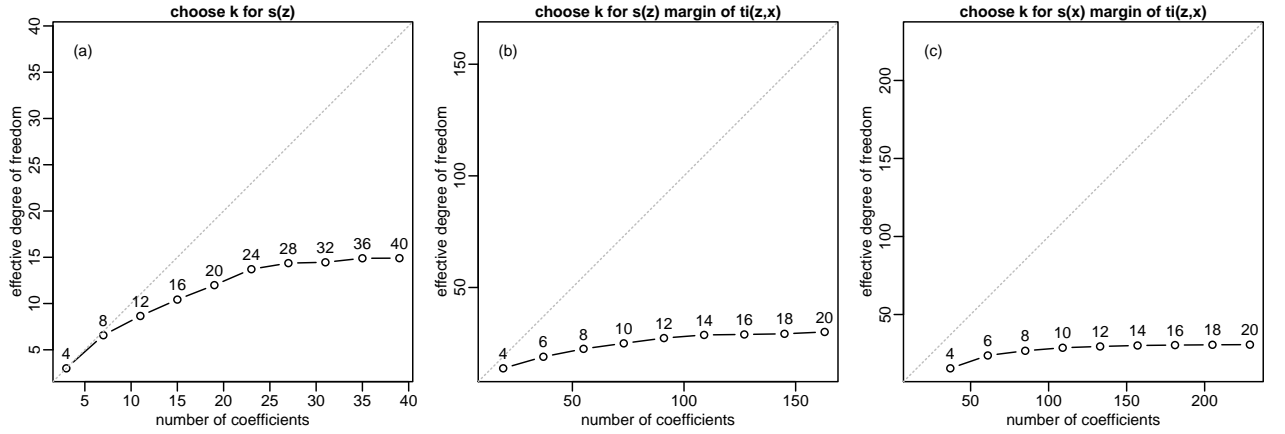
When the model error  $e_i$  is an AR(1) process

$$e_i = \omega \cdot e_{i-1} + \epsilon_i, \quad \epsilon_i \sim N(0, (1 - \omega^2)\phi), \quad i.i.d.$$

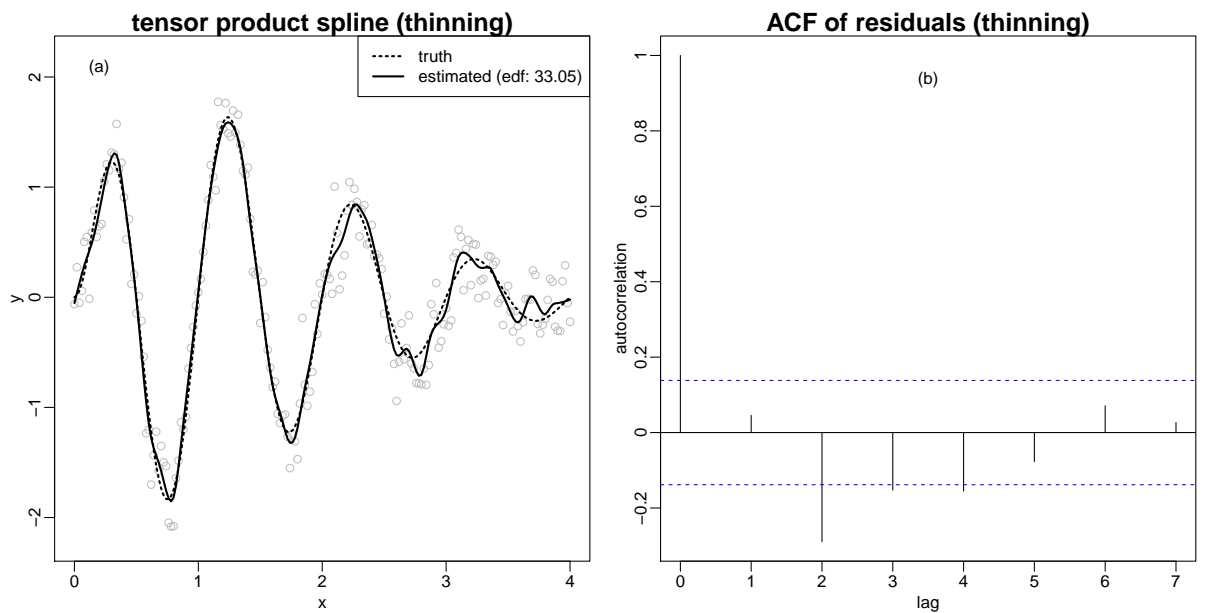




**Figure 2.9:** Choosing  $k$  can be misleading when working on data with autocorrelation. The degree of freedom (*edf*) may never plateau (see panel (b)).



**Figure 2.10:**  $k$  can be properly chosen using thinned data indexed by  $i = 1, 5, 9, \dots, 801$ .



**Figure 2.11:** Fitting the tensor product model 2.3 for thinned data points indexed by  $i = 1, 5, 9, \dots, 801$  reduces overfitting. Residuals also do not exhibit autocorrelation. (Compare with Figure 2.8)

the resulting smoothing or additive model (1.3) has a symmetric tri-diagonal weight matrix (Rue and Held, 2005, §1.1.1):

$$\mathbf{W} = \begin{pmatrix} 1 & -\omega & & & \\ -\omega & 1 + \omega^2 & -\omega & & \\ & \ddots & \ddots & \ddots & \\ & & -\omega & 1 + \omega^2 & -\omega \\ & & & -\omega & 1 \end{pmatrix} \frac{1}{1 - \omega^2}, \quad (2.1)$$

with an upper triangular Cholesky factor

$$\mathbf{W}^{\frac{1}{2}} = \begin{pmatrix} 1 & -\omega & & & \\ & 1 & -\omega & & \\ & & \ddots & \ddots & \\ & & & 1 & -\omega \\ & & & & \sqrt{1 - \omega^2} \end{pmatrix} \frac{1}{\sqrt{1 - \omega^2}}.$$

If  $\omega$  is known a priori,  $\mathbf{W}^{\frac{1}{2}}$  involves no unknown quantity hence the resulting smoothing or additive model (1.3) can be estimated using methods introduced in §4.2 and §4.3.

The matrix  $\mathbf{W}^{\frac{1}{2}}$  is practically useful for data transformation.

- If  $\mathbf{y}$  is a vector of AR(1) samples, i.e.,  $\text{var}(\mathbf{y}) = \mathbf{W}^{-1}\phi$ , then  $\mathbf{z} = \mathbf{W}^{\frac{1}{2}}\mathbf{y}$  is a vector of i.i.d. samples since  $\text{var}(\mathbf{z}) = \mathbf{I}\phi$ .
- If  $\mathbf{z}$  is a vector of i.i.d. samples with  $\text{var}(\mathbf{z}) = \mathbf{I}\phi$ , then  $\mathbf{y} = \mathbf{W}^{-\frac{1}{2}}\mathbf{z}$  is a vector of AR(1) samples since  $\text{var}(\mathbf{y}) = \mathbf{W}^{-1}\phi$ .

The first relation is often used for residual standardization, while the second is often used for simulating AR(1) samples.

First of all, at the beginning of §4.3 it is mentioned that the  $\log |\mathbf{W}|$  in the full REML score (1.9) is dropped, so the REML score reported at convergence of the Newton-Raphson iteration (see Figure 4.9 if you need a revision) needs be “corrected” by adding back  $\log |\mathbf{W}|$  which contains the information of  $\omega$ . This is actually very straightforward. From the previous section, §2.2.2, we can derive  $|\mathbf{W}^{\frac{1}{2}}| = (1 - \omega^2)^{-\frac{n-1}{2}}$ , so  $\log |\mathbf{W}| = 2 \log |\mathbf{W}^{\frac{1}{2}}| = -(n - 1) \log(1 - \omega^2)$ . In the following I will denote the corrected REML score by  $\tilde{\mathcal{V}}_r(\omega)$ .

### 2.2.3 Golden-section search for point estimation of $\omega$

Generally  $\omega$  is unknown and needs be estimated as part of model estimation. An obvious solution is to perform a grid search on  $\omega$ . For example, we fit the smoothing or additive model (1.3) for each of  $\omega = 0.01, 0.02, 0.03, \dots, 0.99$ , then choose the one that minimizes the REML score. However, the grid search is very inefficient. For example, if we search from 0.01 to 0.99, the interval that contains the minimizer is  $[0.01, 0.99], [0.02, 0.99], [0.03, 0.99], \dots$  which shrinks toward one direction slowly after each trial. A better idea is to use something like a bisection that is able to squeeze the interval from both directions. Suppose we evaluate  $\tilde{\mathcal{V}}_r(\omega)$  at 0.25 and 0.75, then we can exclude one of the three intervals  $[0.01, 0.25], [0.25, 0.75]$  and  $[0.75, 0.99]$  from further searching, based on the relationship between  $\tilde{\mathcal{V}}_r(0.25)$  and  $\tilde{\mathcal{V}}_r(0.75)$ .

1. If  $\tilde{\mathcal{V}}_r(0.25) > \tilde{\mathcal{V}}_r(0.75)$ , the minimizer can not be within  $[0.01, 0.25]$ , otherwise  $\tilde{\mathcal{V}}_r(\omega)$  should be non-decreasing on  $[0.25, 0.99]$ , violating  $\tilde{\mathcal{V}}_r(0.25) > \tilde{\mathcal{V}}_r(0.75)$ . Therefore further searching can be restricted on  $[0.25, 0.99]$ .

---

```

 $r = (\sqrt{5} - 1)/2$ 
initialization
 $c = b - r(b - a)$ 
 $d = a + r(b - a)$ 
 $f_c = f(c)$ 
 $f_d = f(d)$ 
start iteration
while ( $|c - d| > \epsilon$ )
    if ( $f_c > f_d$ )
        shrinkage  $[a, b]$  to  $[c, b]$ 
         $a = c$ 
         $c = d$ 
         $f_c = f_d$ 
         $d = a + r(b - a)$ 
         $f_d = f(d)$ 
    else
        shrinkage  $[a, b]$  to  $[a, d]$ 
         $b = d$ 
         $d = c$ 
         $f_d = f_c$ 
         $c = b - r(b - a)$ 
         $f_c = f(c)$ 
return the mid point of  $[a, b]$  as the minimizer
 $(a + b)/2$ 

```

---

**Figure 2.12:** Golden-section search algorithm for minimizing a univariate objective function  $f(x)$  on interval  $[a, b]$ . The minimizer is always inside  $[a, b]$  on each update of this searching interval. In the end the mid point of this interval is returned as the minimizer.  $\epsilon$  is a parameter that controls the precision of result.

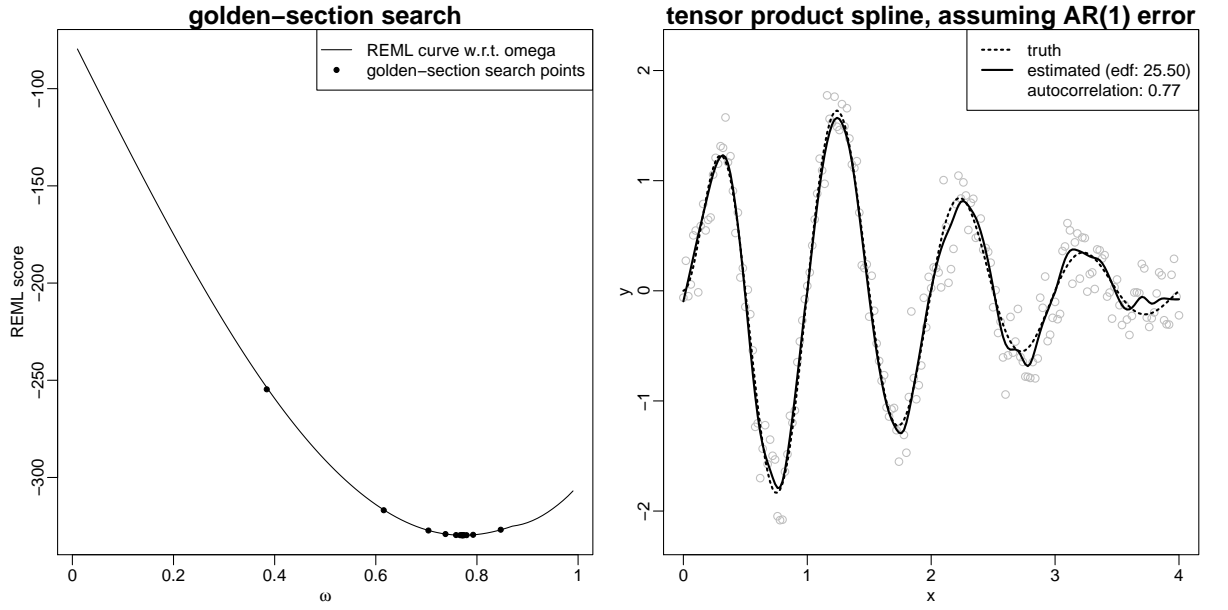
2. If  $\tilde{\mathcal{V}}_r(0.25) < \tilde{\mathcal{V}}_r(0.75)$ , the minimizer can not be within  $[0.75, 0.99]$ , otherwise  $\tilde{\mathcal{V}}_r(\omega)$  should be non-increasing on  $[0.01, 0.75]$ , violating  $\tilde{\mathcal{V}}_r(0.25) < \tilde{\mathcal{V}}_r(0.75)$ . Therefore further searching can be restricted on  $[0.01, 0.75]$ .

This interval shrinkage goes on iteratively. In general, we take two different points  $c$  and  $d$  in interval  $[a, b]$ , then squeeze the new searching interval to either  $[c, b]$  or  $[a, d]$  based on the relationship between  $\tilde{\mathcal{V}}_r(c)$  and  $\tilde{\mathcal{V}}_r(d)$ . There are many ways to choose  $c$  and  $d$  (for example, we can set them randomly), but there exists an optimal choice (in terms of convergence speed) for  $c$  and  $d$ , given by  $c = b - r(b - a)$  and  $d = a + r(b - a)$ , where  $r = (\sqrt{5} - 1)/2$ . This is known as *golden-section search* since  $1/r$  is the golden ratio. Figure 2.12 is a snippet of this algorithm for finding the minimizer of a univariate function  $f$  on  $[a, b]$ .

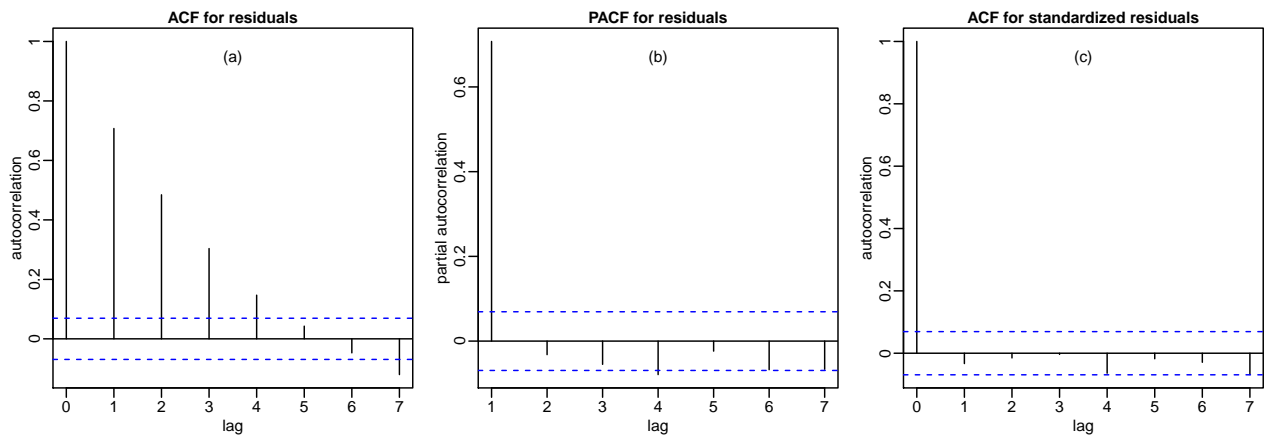
We can now apply the golden-section search to integrate estimation of  $\omega$  in our estimation of the tensor product model 2.3 for the simulated dataset in this section (§2.2). Figure 2.13 illustrates this, and the method proves successful. The point estimate  $\hat{\omega} \approx 0.77$  is close to true value 0.75, and the fitted spline does not exhibit evident overfitting. Figure 2.14 further give the ACF and PACF of raw model residuals of the fitted model, as well as the ACF or standardized residuals. It is clear that the fitted model is adequate since there is no unmodelled autocorrelation in standardized residuals.

## 2.3 Visualization of a tensor product spline

Readers might be interested in knowing what the tensor product spline  $f_1(z, x)$  in the fitted model 2.3 looks like. Figure 2.15 visualizes this surface function with a perspective plot, where surface facets

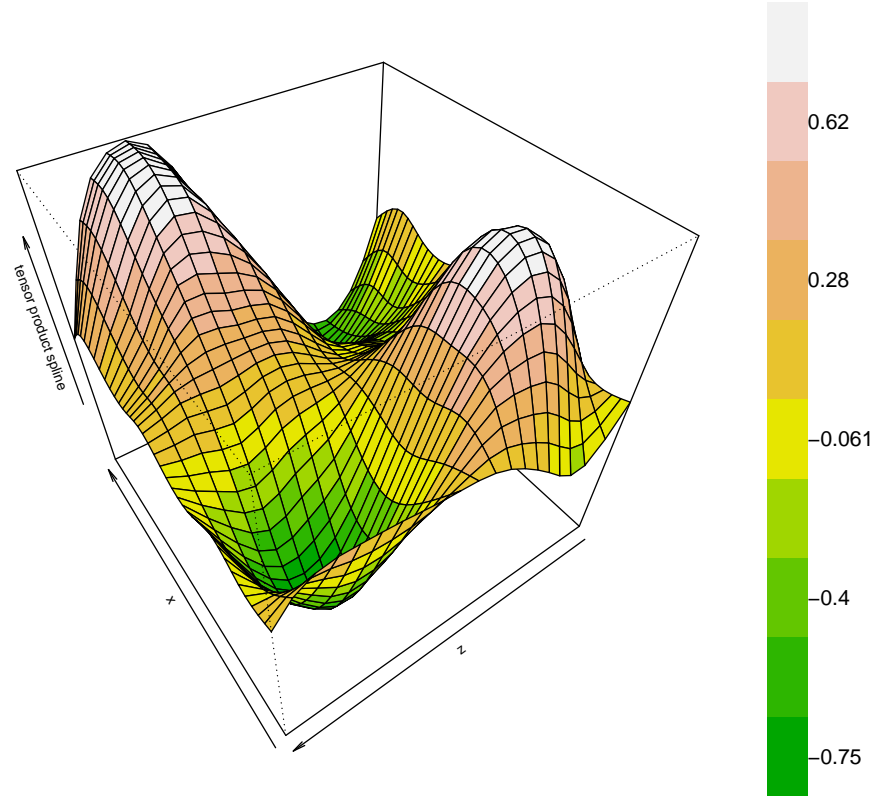


**Figure 2.13:** Implementation of golden-section search to integrate estimation of  $\omega$  in estimation of the tensor product model 2.3 for the dataset in section (§2.2). The left panel sketches the curve of  $\hat{\mathcal{V}}_r(\omega)$ , with the search path of golden-section marked as solid dots on the curve. Convergence occurs after 14 iterations, giving a point estimate of  $\hat{\omega} = 0.7711866$  which is very close to the true value 0.75. The right panel illustrates the fitted spline function. The overfitting problem previously seen in Figure 2.8 is much alleviated.



**Figure 2.14:** Residual inspection of fitted model 2.3 (with golden-section search integrated) for the dataset in section (§2.2). Panel (a) and (b) are respectively ACF and PACF of raw residuals, which convince the AR(1) correlation. Panel (c) is ACF of the standardized residuals (see §2.2.2 for how standardization is performed). The fitted model is adequate as there is no unmodelled autocorrelation.

are coloured according to function values. In the rest of this thesis, this will be the visualization method for any 2D surface smooth functions.

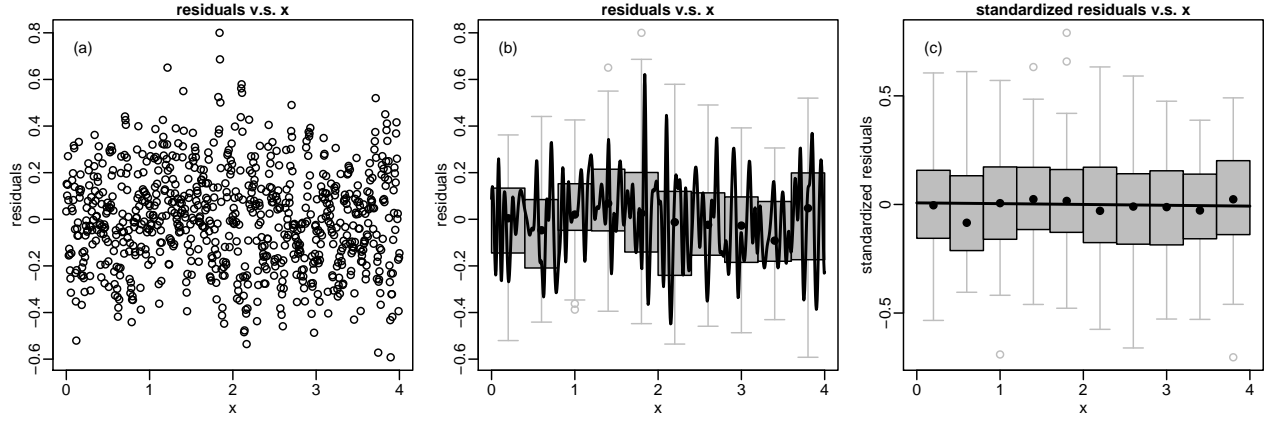


**Figure 2.15:** Visualization of the tensor product spline  $f_1(z, x)$  in the fitted model 2.3 with a perspective plot, where surface facets are coloured according to function values (see the colour bar on the right for reference).

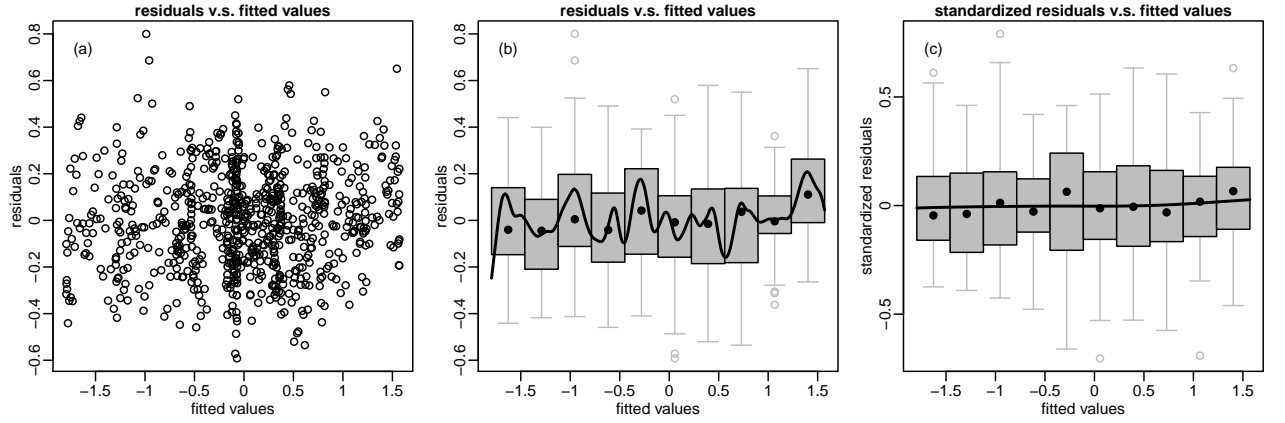
## 2.4 Model checking

In this section I will perform some model checking on the fitted time series model in §2.2.

First of all, let us check whether there is unmodelled trend w.r.t. a variable, say  $x$ . This is commonly visualized by plotting residuals against  $x$ , as is done in panel (a) of Figure 2.16. The scatter plot method is generally inconvenient if there are a great many of data points, as producing the plot takes much time and the resulting cloud of points obscures the interpretation. In this case, it is a good idea to bin the data points and produce a plot for summarized data. For example, panel (b) of Figure 2.16 cuts  $x$  values into 10 bins; residuals are also grouped by bins. Then 25% and 75% sample quantiles for residuals are computed on each bin. A gray shaded box whose height extends to those quantiles is displayed. Sample median is marked by a bold black dot in a box. The whiskers add 1.5 times the interquartile (IQR) to the 75% sample quantile (the top of a box) and subtract 1.5 times the IQR from the 25% sample quantile (the bottom of a box). In the graph, samples within whiskers are not displayed at all; those outside the whiskers are seen as outliers and shown by gray dots. This plotting method substantially reduces the number of points to display, yet well illustrating the variability of data. A smoothing spline is fitted to the original data values (not summarized ones) and overlaid on the plot, for evidence of any trend. From the graph it appears that there is very complicated trend, however, this is because that the smoothing spline is smoothing the AR(1) residuals. Panel (c) produces the same plot using standardized residuals, and it is now clear that there is no unmodelled trend conditional on the AR(1) correlation.



**Figure 2.16:** Inspecting residuals of fitted model 2.3 (with golden-section search integrated) for example data in §2.2 for any unmodelled trend w.r.t. variable  $x$ . Panel (a) is a straightforward scatter plot, while the boxplot-alike visualization in panel (b) and panel (c) are more useful for large dataset. See main text for detailed explanation on those plots. The over complicated trend spotted in panel (b) is due to the fact that the smoothing spline is fitted the AR(1) process in residuals. From panel (c), it seems that there is no unmodelled trend conditional on the AR(1) correlation.



**Figure 2.17:** Inspecting residuals of fitted model 2.3 (with golden-section search integrated) for example data in §2.2 for constant error variance. Panel (a) is a straightforward scatter plot, while the boxplot-alike visualization in panel (b) and panel (c) are more useful for large dataset. See main text for detailed explanation on those plots. The spurious dependence of variance of mean spotted in panel (b) is due to the fact that the smoothing spline is fitted the AR(1) process in residuals. From panel (c), it seems that there is no violation of mean-variance independence conditional on the AR(1) correlation.

Secondly, let us check that residuals have constant variance (i.e., variance is independent of mean). This is often visualized by sketching residuals against fitted values. Panel (a) of Figure 2.17 is a straightforward scatter plot, while panel (b) and panel (c) are respectively boxplot-alike visualization using residuals and standardized residuals. Again, there does not seem to be any dependence of variance on mean.

## 2.5 Summary

In this Chapter we get some hands-on experience on smoothing and additive models via penalized regression splines.

In §2.1.1 I have demonstrated how to model time-varying cyclic pattern in data using a tensor product spline between a cubic cyclic spline and cubic regression spline. This model representation idea can be very useful in modelling time series data with trend and seasonality. While it might first appear that a simpler model representation via a single cubic regression spline is more straightforward, it is

in fact more sensitive to sampling locations of data hence less robust to unevenly spaced sampling, as is illustrated in §2.1.2.

The process of finding appropriate number of knots,  $k$ , for all splines in a model as is demonstrated in §2.1.3 is tedious, and may in fact be time consuming. This contributes to the non-negligible hidden costs for nonparametric modelling in practice.

The data thinning trick introduced in §2.2.1 is very useful for exploratory analysis of data with autocorrelation. In particular, when choosing  $k$  for splines in a model, we are less likely to overestimate the model complexity using thinned data. Correctly assessing model complexity is very important for model development. If an overcomplicated initial model is built at the beginning, the computational burden will be carried forward when new model terms are added and tested.

The Golden-section search method and implementation presented in §2.2.3 integrates the estimation of splines and the estimation of autocorrelation coefficient  $\omega$ . Let  $\mathcal{V}_r(\boldsymbol{\theta})$  be the REML score used in §4.3, the full REML score (in line with (1.9)) is then  $\tilde{\mathcal{V}}_r(\boldsymbol{\theta}, \omega) = \mathcal{V}_r(\boldsymbol{\theta}) - \log |\mathbf{W}(\omega)|$ . The implementation minimizes  $\tilde{\mathcal{V}}_r(\boldsymbol{\theta}, \omega)$  w.r.t.  $(\boldsymbol{\theta}, \omega)$  by a nested iteration. The outer iteration is a one-dimensional minimization w.r.t.  $\omega$ . Then for any trial value  $\omega^{[k]}$ , the inner iteration is a  $q$ -dimensional minimization w.r.t.  $\boldsymbol{\theta}$  (assuming that the vector  $\boldsymbol{\theta}$  has  $q$  elements):

$$\boldsymbol{\theta}^{[k]}(\omega^{[k]}) = \arg_{\boldsymbol{\theta}} \min \{ \mathcal{V}_r(\boldsymbol{\theta}) - \log |\mathbf{W}(\omega^{[k]})| \} = \arg_{\boldsymbol{\theta}} \min \mathcal{V}_r(\boldsymbol{\theta}).$$

Given this  $\boldsymbol{\theta}^{[k]}(\omega^{[k]})$ , the full REML score becomes a univariate objective function:

$$\tilde{\mathcal{V}}_r(\boldsymbol{\theta}^{[k]}(\omega^{[k]}), \omega^{[k]}) = \mathcal{V}_r(\boldsymbol{\theta}^{[k]}(\omega^{[k]})) - \log |\mathbf{W}(\omega^{[k]})|,$$

so that the outer iteration can proceed. The inner iteration is performed with Newton-Raphson algorithm, and the outer iteration is performed with golden-section search. In principle this is not the best idea for such joint estimation. We may derive derivatives of  $\tilde{\mathcal{V}}_r(\boldsymbol{\theta}, \omega)$  w.r.t.  $(\boldsymbol{\theta}, \omega)$  so that a single  $(q + 1)$ -dimensional minimization can be performed by Newton-Raphson algorithm. Unfortunately, this joint estimation introduces dependency between  $\omega$  and regression coefficients, i.e., we now have  $\hat{\beta}_{\boldsymbol{\lambda}, \omega}$  instead of just  $\hat{\beta}_{\boldsymbol{\lambda}}$ , so that REML derivatives presented in §4.3 are no longer valid. While it is possible to derive all required derivatives, the work is non-trivial. By contrast, the nested iteration implementation is convenient. Firstly, nothing needs be changed to the REML estimation procedure introduced in §4.3; secondly, the outer iteration, the golden-section search, requires no derivative computation at all. The only drawback is that model estimation is very computationally expensive. If the golden-section search takes 14 steps for convergence, we need to perform model matrix reduction for 14 times. And if REML estimation for every trial  $\omega^{[k]}$  takes 16 steps for convergence, derivatives of REML score need be computed  $14 \times 16 = 224$  times in total! I will make further remarks on this issue in later Chapters.

In §2.3 and §2.4 some data visualization methods are described. The boxplot alike plot as in Figure 2.16 and Figure 2.17 is a very useful alternative to a large scatter plot with thousands or even millions of dots (i.e., a “cloud” of dots).

In the next Chapter, I will start investigating the Black Smoke dataset.

## Chapter 3

# Preliminary modelling of log Black Smoke (logBS)

In this Chapter, I will start building additive models for Black Smoke.

Firstly, the daily Black Smoke dataset with a number of variables will be introduced in §3.1.1. As a background, some introduction and analysis of the monitoring network is also provided.

Then in the subsequent sections, I will investigate Black Smoke from a few different angles.

- §3.2 will explore the Black Smoke observed at a single station. This mainly involves model development for time series data. In this regard, the time series example previously studied in §2.2 offers a good guidance. Measurements from station *Manchester 11* is used for a thorough case study, because this site had the longest monitoring history with the greatest number of data. If a reasonable model can be built for this station, there wouldn't be any difficulty adapting the same model to other stations.
- §3.3 will apply the time series model concluded from *Manchester 11* to a representative set of stations then pool the result. This also motivates a joint spatial-temporal modelling for all stations' data.
- §3.4 will introduce spatial-temporal modelling using annual mean Black Smoke. Daily Black Smoke data has high day-to-day variability and strong day-to-day correlation, however, once they are aggregated on a yearly basis, the variability of annual mean is substantially reduced and autocorrelation eliminated. The resulting dataset thus only has spatial autocorrelation, making spatial-temporal modelling less challenging.
- §3.5 will tackle a more difficult spatial-temporal modelling task, by looking at daily Black Smoke in year 1967. This year is chosen because it was when the Network had the most measurements: it had the most operating stations and most of them worked for 365 days. Models are fitted for logBS from each day of week to eliminate the impact of temporal autocorrelation in daily observations. It turns out much harder to adequately model space-time relationship without using odd-looking three-way interactions.

While these preliminary modelling exposes a few statistical questions on modelling Black Smoke, they all give way to the computational hurdles of GAM fitting. Modelling daily logBS requires high model complexity (i.e., great number of parameters) that the existing GAM computational engine is not able to deal with.



## 3.1 Introduction to daily logBS dataset and monitoring network

### 3.1.1 The daily logBS dataset

The raw daily Black Smoke dataset simply has three variables:

- the name of the monitoring station, like *Bath 6*, Manchester 11, etc;
- the date of the measurement, like the first day 1961-10-03 and the last day 2005-12-31;
- the Black Smoke measurements, which are non-negative integers.

Plenty of time in my PhD was spent in organizing and enriching the dataset. I was able to obtain the following key information / variables:

- A table providing spatial locations of all monitoring stations. Locations are given by Ordnance Survey grid reference as Northing and Easting. They are originally in metres but are converted to kilometres in my dataset.
- A polygon shapefile for boundaries of British isles. I have processed and exported it as a simple two-column matrix providing the locations of polygon vertices in Northing and Easting.
- Elevation at all stations, extracted from Ordnance Survey's Terrain 50 model (Terrain-50, 2015).
- Daily temperature variables and monthly rainfall variable, extracted from UKCP09 gridded datasets (Perry et al., 2009; Perry and Hollis, 2006) produced by Met Office.

In my organized Black Smoke dataset, Black Smoke measurements (BS) are transformed to logarithmic scale by:  $\log BS = \log(BS + \delta)$ . BS are non-negative integers which are highly skewed (with mean 47.2 and median 21). A log-transformation would make data more conformable to normality, thus additive models for logBS with Gaussian error assumption becomes appropriate. The choice of  $\delta$  is essentially arbitrary, as long as it is not too big compared with BS. I have chosen  $\delta = 1$  so that logBS has a minimum of 0, with mean 3.13 and median 3.09.

Date has been mapped to a number of time variables. There are two types of time variables. The first type is *cumulative* time variable, denoted by a single lowercase letter. These variables are:

- **d**, day since 1961-10-03. For example, 1961-10-03 is day 1, and 2005-12-31 is day 16131;
- **w**, week since 1961-10-02 (the nearest Monday to 1961-10-03), ranging from 1 to 2305.
- **m**, month since 1961-Jan, up to 2005-Dec. For example, 1961-10-03 is in month 10, 1962-10-03 is in month 22 and 2005-12-01 is in month 540.
- **y**, year, ranging from 1961 to 2005.

The second type is *nested* time variable, denoted by a fashion of  $u_v$ , which means  $u$  of  $v$ . These variables are:

- **d<sub>y</sub>**, day of year, ranging from 1 to 366. For example, 1965-01-01 and 1975-01-01 are both day 1 of a year. 2000-12-31 is day 366 of the leap year 2000, while 2003-12-31 is day 365 of the common year 2003.

- $d_w$ , day of week, ranging from 1 to 7, for Monday to Sunday.
- $w_y$ , week of year, ranging from 0 to 53. This variable needs some explanation. I define that the first week of a year starts from the nearest Monday to the first day of that year. For example, 1989-01-01 is Sunday, so it is in week 0 of the year and week 1 starts from 1989-01-02. As another example, 2019-01-01 is Tuesday, so it is week 1 of the year and there is no week 0.

The dataset has the following variables associated with monitoring stations:

- $i$ , numeric station ID, from 1 to 2874.
- $\{e_i, n_i\}$ , spatial locations of station  $i$ ;
- $h_i$ , elevation of station  $i$ ;
- $E_i$ , environmental type at  $i$ . According to Loader (2002), environmental type has five major levels: residential area with high-density housing (A), suburb residential area with medium-density housing (B), industrial area (C), commercial area or town centre (D) and rural community or open country (R).

For meteorological variables, the dataset has

- $T_{id}^0$ , daily minimum temperature ( $^{\circ}\text{C}$ ) for station  $i$  and day  $d$ ;
- $T_{id}^1$ , daily maximum temperature ( $^{\circ}\text{C}$ ) for station  $i$  and day  $d$ ;
- $T_{id}^*$ , diurnal temperature variation for station  $i$  and day  $d$ , defined by  $T_{id}^1 - T_{id}^0$ ;
- $r_{im}$ , monthly total rainfall (millimetre) for station  $i$  and month  $m$ .

In total, the dataset has 9489903 (about 10 million) data from a total of 2874 stations for 45 years from 1961-10-03 to 2005-12-31.

### 3.1.2 The Black Smoke monitoring network

Although the dataset has a big volume of data, the information over space and time was only sufficient for 20 years prior to 1981 or 1982. Figure 3.1 maps the Network on 1975, 1985, 1995 and 2005. To enhance visualization a gray shaded area within 10 km distance from all 2874 historical stations (I will hereafter call this the “local domain” of the Network) is also displayed. First of all, stations were not evenly spread out over the islands; secondly, even on the local domain the Network was getting thinner and thinner in later years.

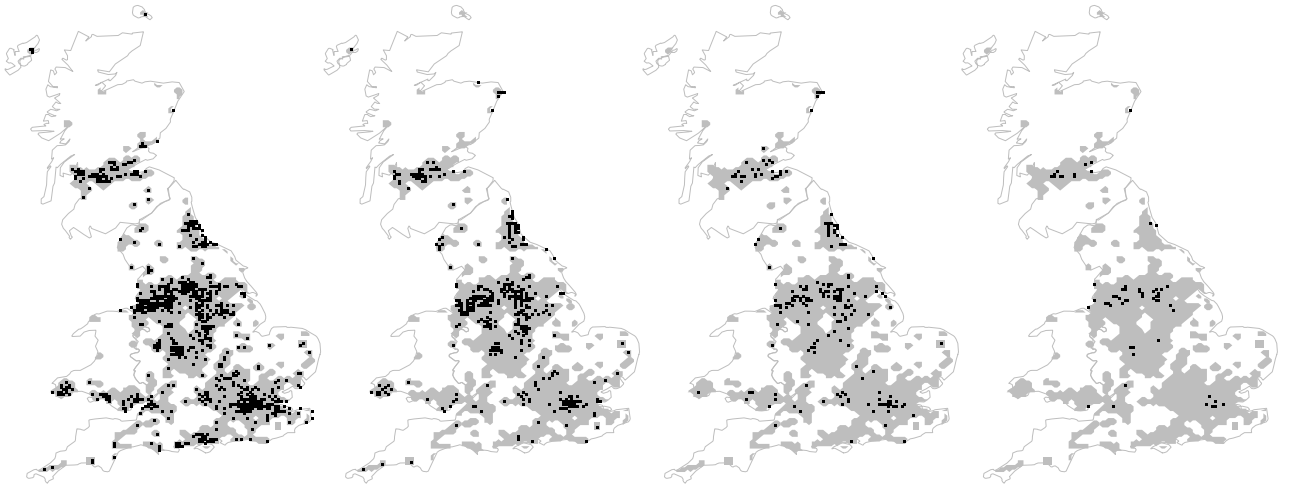
Since Black Smoke declined over decades, It has been suspected that stations with lower Black Smoke observations were more likely to be shut down through years, making the Network sparser and sparser. Under the assumption that logBS declined with the same rate for all stations, then if this preferential sites closure existed, the sample spatial mean or other similar statistics summarized from operating stations alone should decline slower and slower, over even increase at some point. A common statistic for monitoring networks is *exceedance*, the number of days when Black Smoke exceeded a certain limit. Table 3.1 lists the Black Smoke limits imposed by EC Directive 80/779/EEC. This Directive was introduced to UK in 1980, exactly when the Network experienced a massive wave of sites’ closure. The role of the Network had changed thereafter. The Network manual states clearly that monitoring compliance with this Directive became the Network’s main objective throughout the remainder of its

1237 sites in year 1975

563 sites in year 1985

225 sites in year 1995

65 sites in year 2005



**Figure 3.1:** Black Smoke monitoring network (1961-2005). The maps outline the Network in 1975, 1985, 1995 and 2005 by landmarking operational sites in black dots. The gray shaded area covers regions within 10 km distance from all 2874 historical stations (I will call it a “local domain” of the Network).

**Table 3.1:** Directive 80/779/EEC Limits for Black Smoke (1980 ~ 2004)

reference period	Limits ( $\mu\text{gm}^{-3}$ )
Year (median of daily values)	68
Winter (median of daily values Oct. - Mar.)	111
Year (Peak: 98 Percentile of daily values)	213

existence. (In fact, the final closure of the Network was also related to this Directive. On January 1st, 2005, this Directive was repealed, and the Network’s monitoring role was ceased at the end of that year.)

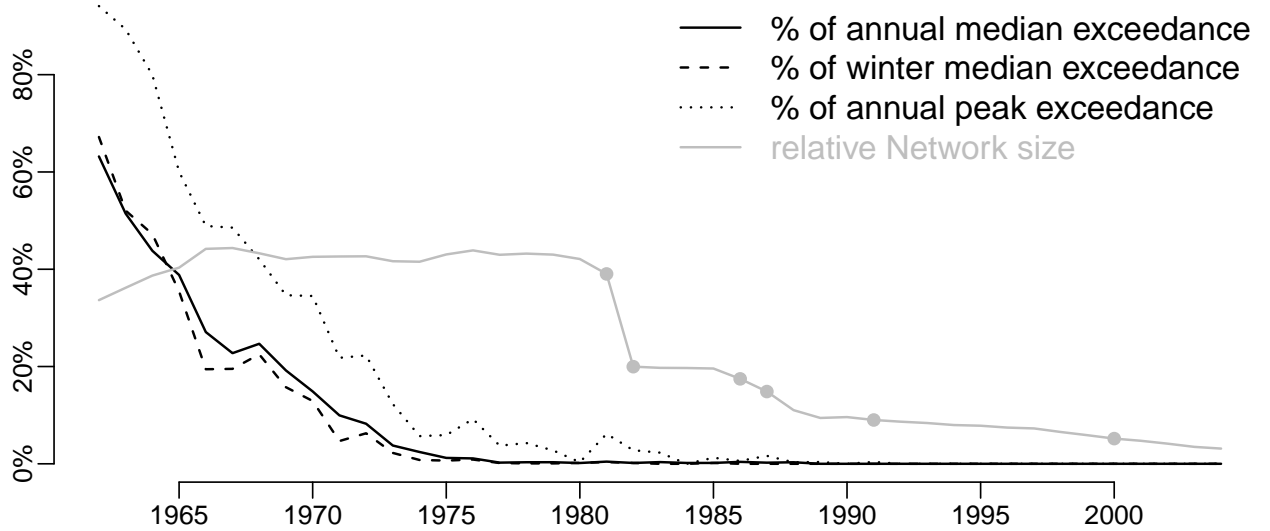
Figure 3.2 sketches the percentage of sites failing to comply the EC Directive limits over years. It is interesting to observe that by 1980 when EC Directive was introduced, all three percentages had declined to almost zero. Then in 1981 and 1982, the percentage associated with annual peak limits suddenly rose up again. Since this was the time when the Network was reduced in size, it is almost surely an evidence that 1) sites where BS had been consistently below the limits for many years were closed; 2) only sites where Black Smoke was not low enough were retained for a few more years for monitoring.

In summary, the information carried by BS observations is unbalanced over the years.

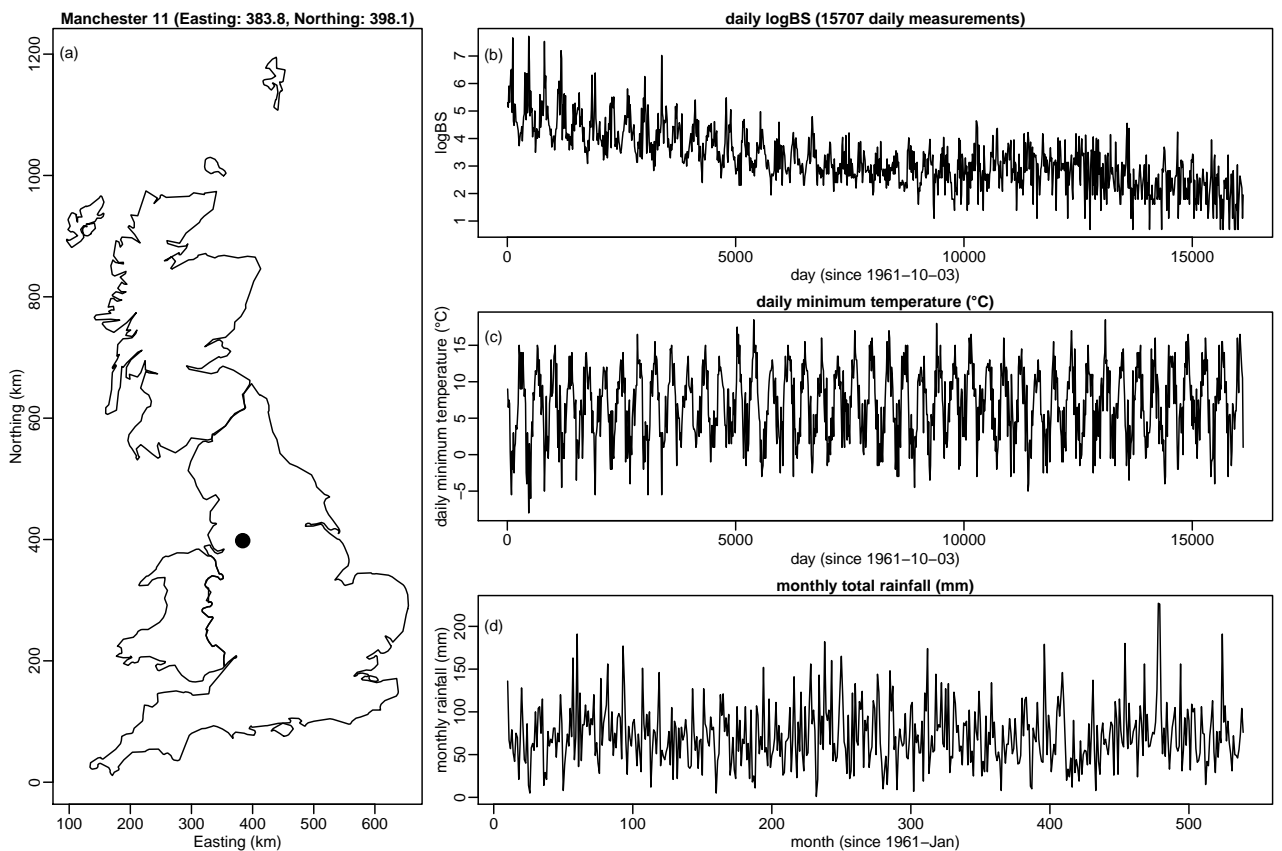
## 3.2 Time series models for daily logBS from “MANCHESTER 11”

### 3.2.1 Description of data

In this section let us build a time series model for logBS from station *Manchester 11*, which had the longest monitoring history for the entire 45 years (1961 ~ 2005) with 15707 data. Figure 3.3 is an illustration of this dataset.



**Figure 3.2:** A profile of the Network. The gray solid line gives the relative size of the Network, defined as the ratio (in percentage) between the number of operational sites in a year and the total number of historical sites. The black solid, dashed and dotted lines respectively give the percentage of sites in a year failing to comply the EC Directive limits in annual median, winter median and annual peak (see Table 3.1).



**Figure 3.3:** Data from monitoring station *Manchester 11*. Panel (a) is the location of this station on map; Panel (b) is the daily logBS time series; Panel (c) is the daily minimum temperature at this location; Panel (d) is the monthly total rainfall at this location.

### 3.2.2 Building a separate model for each day of week

It can be spotted from Panel (a) of Figure 3.3 that the logBS time series has an evident cyclic characteristic within year, i.e., seasonality. Inspired by the case study in §2.1, I propose the following model:

$$\log\text{BS}_{\mathbf{d}} = f_1(\mathbf{d}_y; k_1) + f_2(\mathbf{d}; k_2) + f_3(\mathbf{d}_y, \mathbf{d}; k_{3,1}, k_{3,2}) + \epsilon_{\mathbf{d}},$$

where

- $f_1(\mathbf{d}_y)$  is a cubic cyclic spline with  $k_1$  number of knots, modelling the seasonality of logBS;
- $f_2(\mathbf{d})$  is a natural cubic spline with  $k_2$  number of knots, modelling the trend of logBS;
- $f_3(\mathbf{d}_y, \mathbf{d})$  is a tensor product spline between a cubic cyclic spline margin for  $\mathbf{d}$  with  $k_{3,1}$  knots and a natural cubic spline margin for  $\mathbf{d}$  with  $k_{3,2}$  knots, modelling how seasonality varies with time;
- $\epsilon_{\mathbf{d}}$  is an i.i.d. Gaussian error on each day.

The i.i.d. error assumption is very implausible for daily data that usually have strong autocorrelation. So I will apply the data thinning trick explained in §2.2.1, by looking at data from for example, Monday only. That is, I will start by building a model for each day of week. Note that in this way, it is equivalent to use  $\mathbf{w}_y$  in place of  $\mathbf{d}_y$  and  $\mathbf{w}$  instead of  $\mathbf{d}$ . So, I will write this model as

**Model 3.1:**

$$\log\text{BS}_{\mathbf{w}} = f_1(\mathbf{w}_y; k_1) + f_2(\mathbf{w}; k_2) + f_3(\mathbf{w}_y, \mathbf{w}; k_{3,1}, k_{3,2}) + \epsilon_{\mathbf{w}},$$

While model 3.1 is potentially a reasonable structure on thinned data, its required complexity, i.e., the number of knots needed for each spline is unknown yet. Following the “choose k” procedure explained in §2.1.3, I conducted the following selection process.

1. Fit model 3.1 (for each day of week) with  $k_1 = 10$ ,  $k_2 = 30$ ,  $k_{3,1} = 10$  and  $k_{3,2} = 30$ . Extract the partial residual w.r.t.  $\hat{f}_1(\mathbf{w}_y, 10)$ , then fit model (for each day of week)

$$\text{partial residuals} = f_1(\mathbf{w}_y, k_1) + \text{i.i.d. errors}$$

for a set of increasingly big  $k_1$ . The panel (a) of Figure 3.4 implies that  $k_1 = 15$  is sufficient for all days of week.

2. Fit model 3.1 (for each day of week) with  $k_1 = 15$ ,  $k_2 = 30$ ,  $k_{3,1} = 10$  and  $k_{3,2} = 30$ . Extract the partial residual w.r.t.  $\hat{f}_2(\mathbf{w}, 30)$  then fit model (for each day of week)

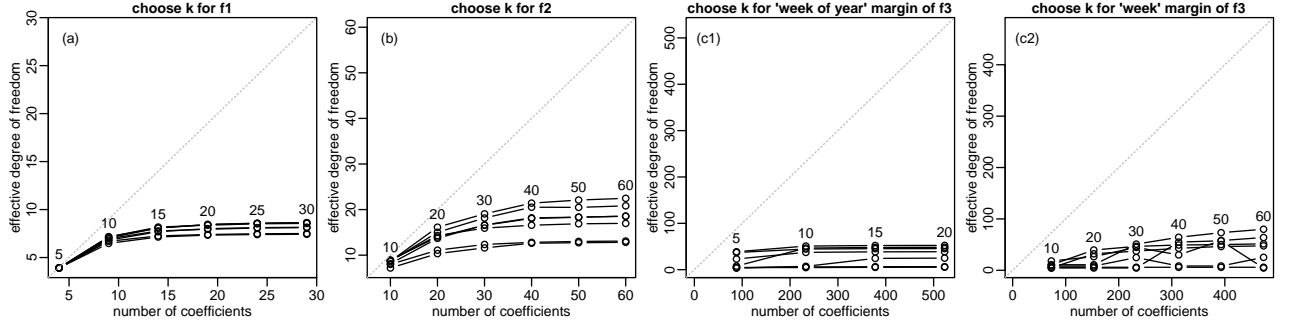
$$\text{partial residuals} = f_2(\mathbf{w}, k_2) + \text{i.i.d. errors}$$

for a set of increasingly big  $k_2$ . The panel (b) of Figure 3.4 implies that  $k_2 = 40$  is sufficient for all days of week.

3. Fit model 3.1 (for each day of week) with  $k_1 = 15$ ,  $k_2 = 40$ ,  $k_{3,1} = 10$  and  $k_{3,2} = 30$ . Extract the partial residual w.r.t.  $\hat{f}_3(\mathbf{w}_y, \mathbf{w}; 10, 30)$  and fit model (for each day of week)

$$\text{partial residuals} = f_3(\mathbf{w}_y, \mathbf{w}; k_{3,1}, 30) + \text{i.i.d. errors}$$

for a set of increasingly big  $k_{3,1}$ . The panel (c1) of Figure 3.4 implies that  $k_{3,1} = 10$  is sufficient for all days of week.



**Figure 3.4:** Choosing  $k_1$ ,  $k_2$ ,  $k_{3d}$  and  $k_{3d}$  for model 3.1. These graphs are similar to those in Figure 2.7, except that the results for all seven days of a week are displayed. This gives a good idea on a reasonable choice of  $k$ -values for all seven models.

**Table 3.2:** Number knots, coefficients and effective degree of freedom in fitted model 3.1, for each day of week.

	No. of knots	No. of coef.	Mon	Tue	Wed	Thu	Fri	Sat	Sun
			effective degree of freedom						
$\hat{f}_1$	$k_1 = 15$	13	7.0	6.1	7.1	6.7	6.3	6.7	7.0
$\hat{f}_2$	$k_2 = 40$	39	15.4	12.3	16.8	20.2	20.0	17.1	11.1
$\hat{f}_3$	$k_3 = (10, 40)$	$8 \times 39 = 312$	62.4	21.0	46.4	41.2	53.3	4.8	47.6

- Update model 3.1 with  $k_1 = 15$ ,  $k_2 = 40$ ,  $k_{3,1} = 10$  and  $k_{3,2} = 30$ . Extract the partial residual w.r.t.  $\hat{f}_3(\mathbf{w}_y, \mathbf{w}; 10, 30)$  and fit model (for each day of week)

$$\text{partial residuals} = f_3(\mathbf{w}_y, \mathbf{w}; 10, k_{3,2}) + \text{i.i.d. errors}$$

for a set of increasingly big  $k_{3,2}$ . The panel (c2) of Figure 3.4 implies that  $k_{3,2} = 40$  is sufficient for all days of week.

So an adequate model is

$$\log \text{BS}_w = f_1(\mathbf{w}_y; 15) + f_2(\mathbf{w}; 40) + f_3(\mathbf{w}_y, \mathbf{w}; 10, 40) + \epsilon_w,$$

and Table 3.2 summarizes the number of parameters and the resulting effective degree of freedom for each splines.

All seven fitted models turn out adequate. Figure 3.5 sketches the sample autocorrelation (ACF) of residuals for all models, and no significant autocorrelation is spotted.

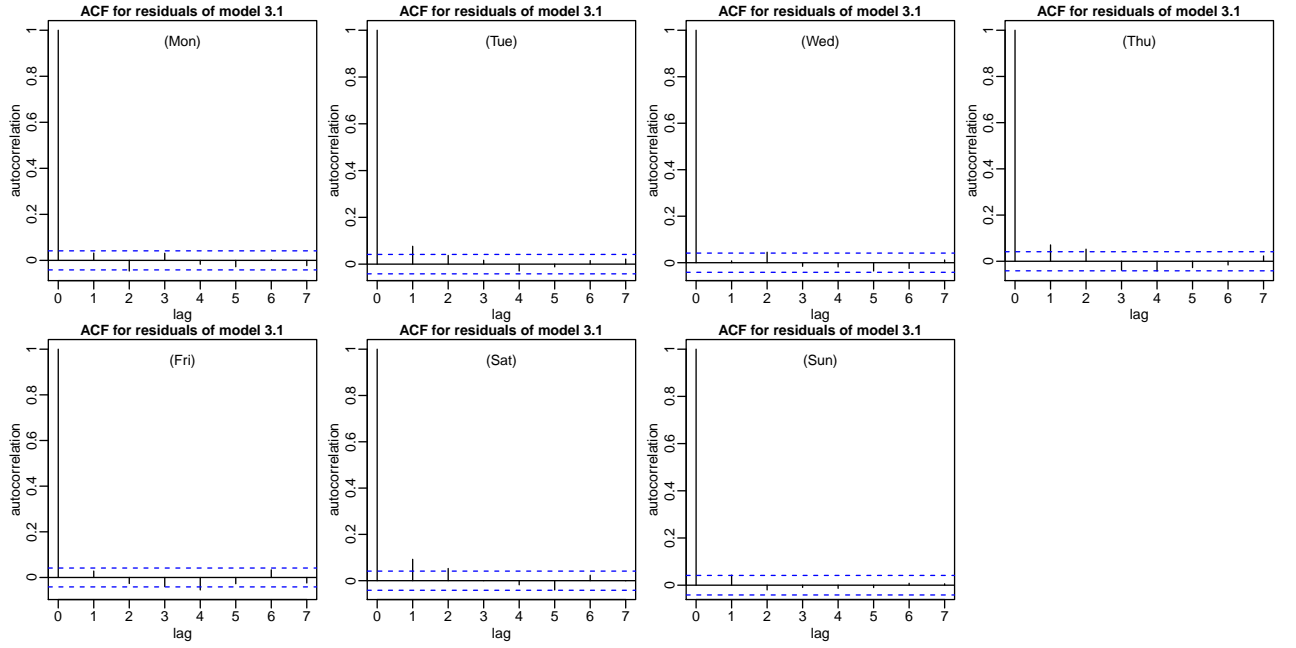
The data thinning trick has indeed helped the selection of  $k$ . If using complete data, the complexity of model 3.1 will be severely overestimated. For example,  $k_1$  must be at least 50 to be seen “sufficient” and  $k_2$  needs be 100. That is, the spline function is fitting the autocorrelation in daily data.

### 3.2.3 A joint model for all days of week

Now I will build a joint model for all days of week. To start with, I introduce  $\mathbf{d}_w$  as a factor of seven levels so that the intercept and all spline functions in model 3.1 can vary between levels. This model can be expressed as:

$$\log \text{BS}_d = \sum_{s=1}^7 \delta_s^{\mathbf{d}_w} u_s + \sum_{s=1}^7 \delta_s^{\mathbf{d}_w} f_{1,s}(\mathbf{w}_y; 15) + \sum_{s=1}^7 \delta_s^{\mathbf{d}_w} f_{2,s}(\mathbf{w}; 40) + \sum_{s=1}^7 \delta_s^{\mathbf{d}_w} f_{3,s}(\mathbf{w}_y, \mathbf{w}; 10, 40) + e_d,$$

where



**Figure 3.5:** Sample autocorrelation (ACF) of residuals for model 3.1 fitted to each day of week. No significant autocorrelation is seen.

- $u_s$  is the intercept for day  $s$  of week (so that there will be  $u_1, u_2, \dots, u_7$ ) and  $\delta_s^{\mathbf{d}_w}$  is a dummy variable with value 1 when  $s = \mathbf{d}_w$  and 0 otherwise;
- $f_{1,s}$  is a cubic cyclic spline of  $\mathbf{w}_y$ , modelling the seasonality of logBS from level  $s$ ;
- $f_{2,s}$  is a natural cubic spline of  $\mathbf{w}$ , modelling the trend of logBS from level  $s$ ;
- and  $f_{3,s}$  is tensor product spline, modelling how seasonality of logBS from level  $s$  varies with time;
- $e_d$  is a model error with some autocorrelation structure.

In the above specification, there is no contrasts applied to  $\mathbf{d}_w$ . An alternative specification is to set Monday as the reference level, so that spline functions for Monday are the reference smooth functions, and spline functions for other levels are the deviation from those reference smooth functions. This gives

### Model 3.2a:

$$\begin{aligned} \log\text{BS}_d = & u_1 + \sum_{s=2}^7 \delta_s^{\mathbf{d}_w} u_s + f_1(\mathbf{w}_y; 15) + \sum_{s=2}^7 \delta_s^{\mathbf{d}_w} f_{1,s}(\mathbf{w}_y; 15) + \\ & f_2(\mathbf{w}; 40) + \sum_{s=2}^7 \delta_s^{\mathbf{d}_w} f_{2,s}(\mathbf{w}; 40) + \\ & f_3(\mathbf{w}_y, \mathbf{w}; 10, 40) + \sum_{s=2}^7 \delta_s^{\mathbf{d}_w} f_{3,s}(\mathbf{w}_y, \mathbf{w}; 10, 40) + e_d, \end{aligned}$$

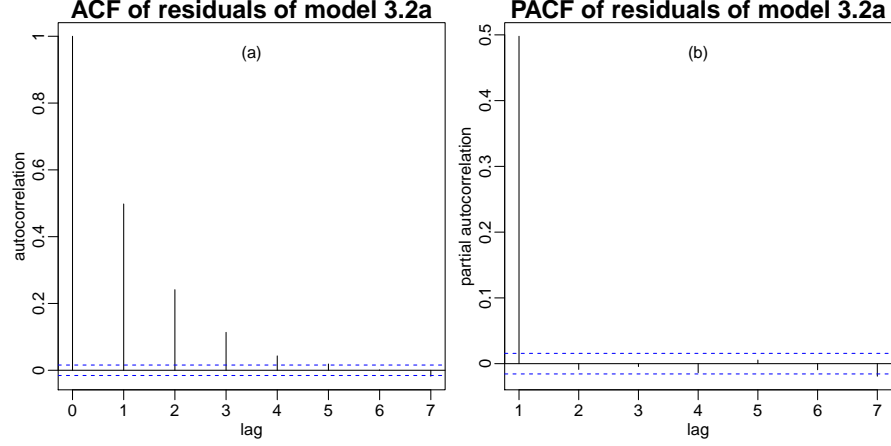
The advantage of this specification in GAM estimation, is that if there is little difference between day  $s$  and Monday, the “difference smooth” will be “penalized out”, ending up with zero degree of freedom. This will be a sign that some components can be dropped from the model.

Fitting this additive model assuming that  $e_d$  is i.i.d. ends up with the effective degree of freedom for all functions as is summarized in Table 3.3. There are the following conclusion.

1.  $f_{3,2}, f_{3,3}, \dots, f_{3,7}$  are “penalized out”, so they can be dropped from the model.

**Table 3.3:** Effective degree of freedom for all functions in fitted model 3.2a, assuming i.i.d. model error.

$\hat{f}_1$	$\hat{f}_{1,2}$	$\hat{f}_{1,3}$	$\hat{f}_{1,4}$	$\hat{f}_{1,5}$	$\hat{f}_{1,6}$	$\hat{f}_{1,7}$
11.58	0.00	0.10	0.75	1.58	0.98	2.23
$\hat{f}_2$	$\hat{f}_{2,2}$	$\hat{f}_{2,3}$	$\hat{f}_{2,4}$	$\hat{f}_{2,5}$	$\hat{f}_{2,6}$	$\hat{f}_{2,7}$
25.33	1.00	2.53	2.90	2.22	1.00	6.58
$\hat{f}_3$	$\hat{f}_{3,2}$	$\hat{f}_{3,3}$	$\hat{f}_{3,4}$	$\hat{f}_{3,5}$	$\hat{f}_{3,6}$	$\hat{f}_{3,7}$
235.79	0.00	0.00	0.00	0.00	0.00	0.00



**Figure 3.6:** Strong residual autocorrelation is spotted in residuals of fitted model 3.2a. Panel (a) is the ACF and panel (b) is the PACF. The cutoff of partial autocorrelation at lag 1 implies that residuals are AR(1) correlated.

2. Not all of  $f_{2,2}, f_{2,3}, \dots, f_{2,7}$  are “penalized out”, so they should all be retained.
3. Compared with Table 3.2, the effective degree of freedom for  $f_1, f_2$  and  $f_3$  have increased. This is a message that those spline functions are fitting autocorrelation in daily data.

In fact, the existence of strong autocorrelation can be verified from model residuals, as is shown in Panel (a) of Figure 3.6. The exponential decay of autocorrelation w.r.t. time lag probably implies that residuals are AR(1) correlated. Panel (b) confirms this with the cutoff of partial autocorrelation at lag 1.

These observations lead to the following model:

**Model 3.2b:**

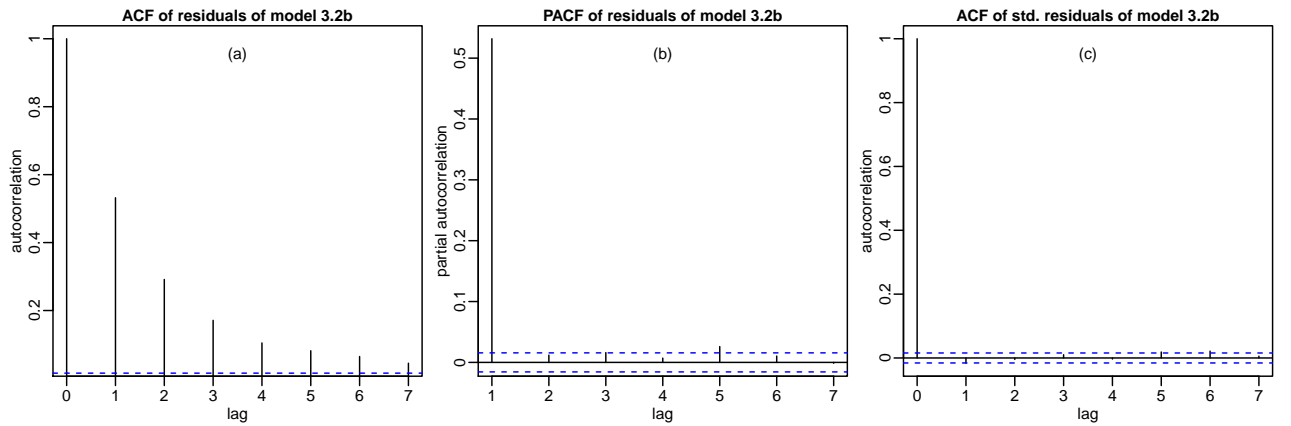
$$\log \text{BS}_d = u_1 + \sum_{s=2}^7 \delta_s^{\text{d}_w} u_s + f_1(\mathbf{w}_y; 15) + \sum_{s=2}^7 \delta_s^{\text{d}_w} f_{1,s}(\mathbf{w}_y; 15) + f_2(\mathbf{w}; 40) + \sum_{s=2}^7 \delta_s^{\text{d}_w} f_{2,s}(\mathbf{w}; 40) + f_3(\mathbf{w}_y, \mathbf{w}; 10, 40) + e_d,$$

where  $f_{3,2}$  to  $f_{3,7}$  in model 3.2a are dropped and  $e_d$  is assumed an AR(1) stochastic process with some unknown autocorrelation coefficient  $\omega$ . The golden-section search implementation introduced in §2.2.3 makes estimation of this additive model possible. It gives a point estimate  $\hat{\omega} = 0.5447645$  (see 3.8 for the search path of golden-section), and the resulting degree of freedom for all spline functions in the model are reported in Table 3.4. Note that the values for  $f_1$  and  $f_2$  are now similar to those in Table 3.2. This is a message that autocorrelation has been properly modelled and the spline functions are not overfitting data. Figure 3.7 is a formal residual checking, where AR(1) assumption is validated.

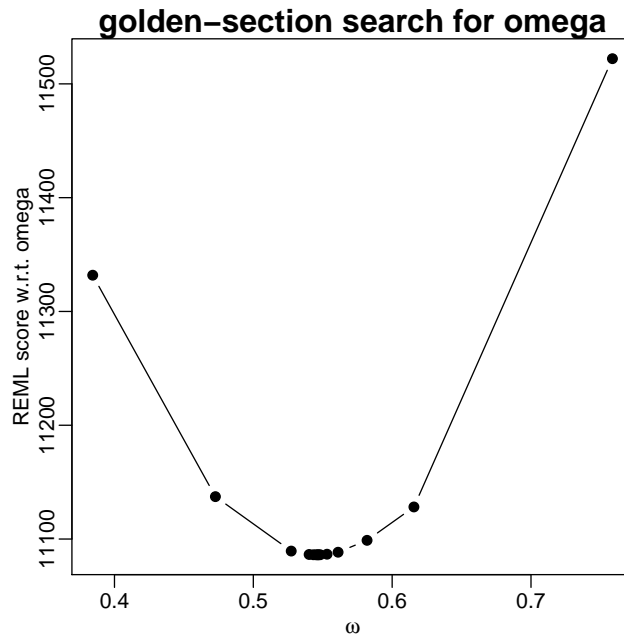


**Table 3.4:** Effective degree of freedom for all functions in fitted model 3.2b, assuming AR(1) error with unknown autocorrelation coefficient.

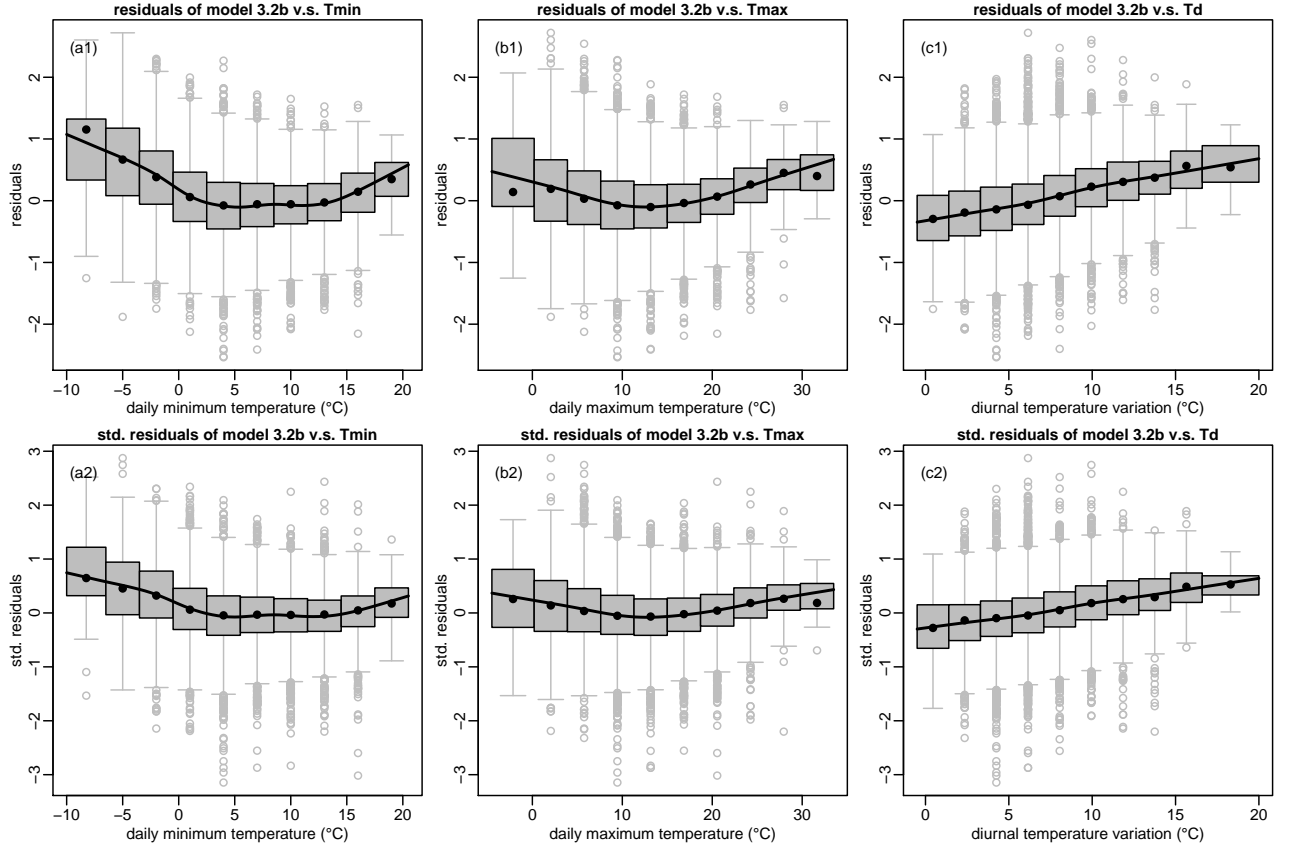
$\hat{f}_1$	$\hat{f}_{1,2}$	$\hat{f}_{1,3}$	$\hat{f}_{1,4}$	$\hat{f}_{1,5}$	$\hat{f}_{1,6}$	$\hat{f}_{1,7}$
7.83	0.00	0.00	0.00	1.40	1.39	2.38
$\hat{f}_2$	$\hat{f}_{2,2}$	$\hat{f}_{2,3}$	$\hat{f}_{2,4}$	$\hat{f}_{2,5}$	$\hat{f}_{2,6}$	$\hat{f}_{2,7}$
19.8	2.25	3.03	3.49	2.75	1.00	8.79
$\hat{f}_3$						
81.65						



**Figure 3.7:** The AR(1) error assumption for model 3.2b is adequate. Although raw residuals are still AR(1) correlated (see panel (a) and panel (b)), standardized residuals are now i.i.d. (see panel (c)). See §2.2.2 if you need a revision on residual standardization.



**Figure 3.8:** The search path of golden-section iterations for fitted model 3.2b. A point estimate of  $\hat{\omega} = 0.5447645$  is obtained after 14 iterations.



**Figure 3.9:** Inspecting residuals of fitted model 3.2b against daily minimum temperature ( $T_{\text{min}}$ )  $T_{\text{d}}^0$ , daily maximum temperature ( $T_{\text{max}}$ )  $T_{\text{d}}^1$  and diurnal temperature variation ( $T_{\text{d}}$ )  $T_{\text{d}}^*$ . Unmodelled trend has been detected so these covariates could improve our model. Since  $T_{\text{d}}^0$  and  $T_{\text{d}}^1$  are highly correlated, they should not enter the model in an additive manner; they need be jointly modelled as a bivariate thin-plate spline. An alternative idea may be including  $T_{\text{d}}^0$  and  $T_{\text{d}}^*$  in an additive manner as these two variables are uncorrelated. The plotting method used is explained in §2.4.

### 3.2.4 Including meteorological covariates

So far no meteorological covariates have been used for model building. It is time to see if they can improve our model. For this *Manchester 11* case study, I will simply denote temperature variables by for example,  $T_{\text{d}}^0$  than  $T_{\text{id}}^0$ , because there is only a single station.

Figure 3.9 sketches residuals and standardized residuals of fitted model 3.2b against daily minimum temperature  $T_{\text{d}}^0$ , daily maximum temperature  $T_{\text{d}}^1$  and diurnal temperature variation  $T_{\text{d}}^*$ . A clear, similar non-linear relationship is seen between residuals and  $T_{\text{d}}^0$  or  $T_{\text{d}}^1$ . This is probably because that  $T_{\text{d}}^0$  and  $T_{\text{d}}^1$  are highly correlated (their Pearson correlation coefficient is 0.82 at *Manchester 11*). This means that if both variables would enter the model, they should not be included in an additive manner. A bivariate thin-plate spline is more appropriate. One motivation of this bivariate function is that interaction between  $T_{\text{d}}^0$  or  $T_{\text{d}}^1$  can also be modelled. Such interaction may make sense. Suppose that two days have the same daily minimum temperature at  $-5^\circ\text{C}$ , but one day has a daily maximum temperature at  $10^\circ\text{C}$  while the other has that at  $0^\circ\text{C}$ . Obviously the second day is colder so there is probably more heating supply, more fuel combustion hence more Smoke emission. In other words,  $T_{\text{d}}^0$  alone may not be fully explain the temperature effect on Smoke concentration; diurnal temperature variation could matter as well. Panel (c1) and panel (c2) of Figure 3.9 seem to verify this. In fact,  $T_{\text{d}}^0$  and  $T_{\text{d}}^*$  can be included in a model in an additive manner because they are basically uncorrelated (their Pearson correlation coefficient is 0.02 at *Manchester 11*). I thus consider a model

**Table 3.5:** Effective degree of freedom for all functions in fitted model 3.3a, assuming AR(1) error with unknown autocorrelation coefficient.

$\hat{f}_1$	$\hat{f}_{1,2}$	$\hat{f}_{1,3}$	$\hat{f}_{1,4}$	$\hat{f}_{1,5}$	$\hat{f}_{1,6}$	$\hat{f}_{1,7}$
9.50	0.00	0.00	0.00	1.14	1.33	2.12
$\hat{f}_2$	$\hat{f}_{2,2}$	$\hat{f}_{2,3}$	$\hat{f}_{2,4}$	$\hat{f}_{2,5}$	$\hat{f}_{2,6}$	$\hat{f}_{2,7}$
20.80	1.79	2.89	3.77	2.45	1.00	8.43
$\hat{f}_3$	$\hat{f}_4$	$\hat{f}_5$				
87.46	8.57	3.74				

**Model 3.3a:** model 3.2b +  $f_4$  +  $f_5$

$$\begin{aligned} \log \text{BS}_d = & u_1 + \sum_{s=2}^7 \delta_s^{\text{d}_w} u_s + f_1(\mathbf{w}_y; 15) + \sum_{s=2}^7 \delta_s^{\text{d}_w} f_{1,s}(\mathbf{w}_y; 15) + \\ & f_2(\mathbf{w}; 40) + \sum_{s=2}^7 \delta_s^{\text{d}_w} f_{2,s}(\mathbf{w}; 40) + f_3(\mathbf{w}_y, \mathbf{w}; 10, 40) + \\ & f_4(\mathbf{T}_d^0; 15) + f_5(\mathbf{T}_d^*; 10) + e_d, \end{aligned}$$

where new components to model 3.2b are

1.  $f_4$  is a natural cubic spline of  $\mathbf{T}_d^0$ , with 15 knots (this seems more than sufficient from Figure 3.9);
2.  $f_5$  is a natural cubic spline of  $\mathbf{T}_d^*$ , with 10 knots (this seems more than sufficient from Figure 3.9; note that while the relationship appears linear from Figure 3.9, it might still be a good idea to model it as a spline, as such relationship can be non-linear for other stations than *Manchester 11* and we should not overly restrict the form representation at this stage);
3.  $e_d$  is an AR(1) model error with unknown autocorrelation coefficient.

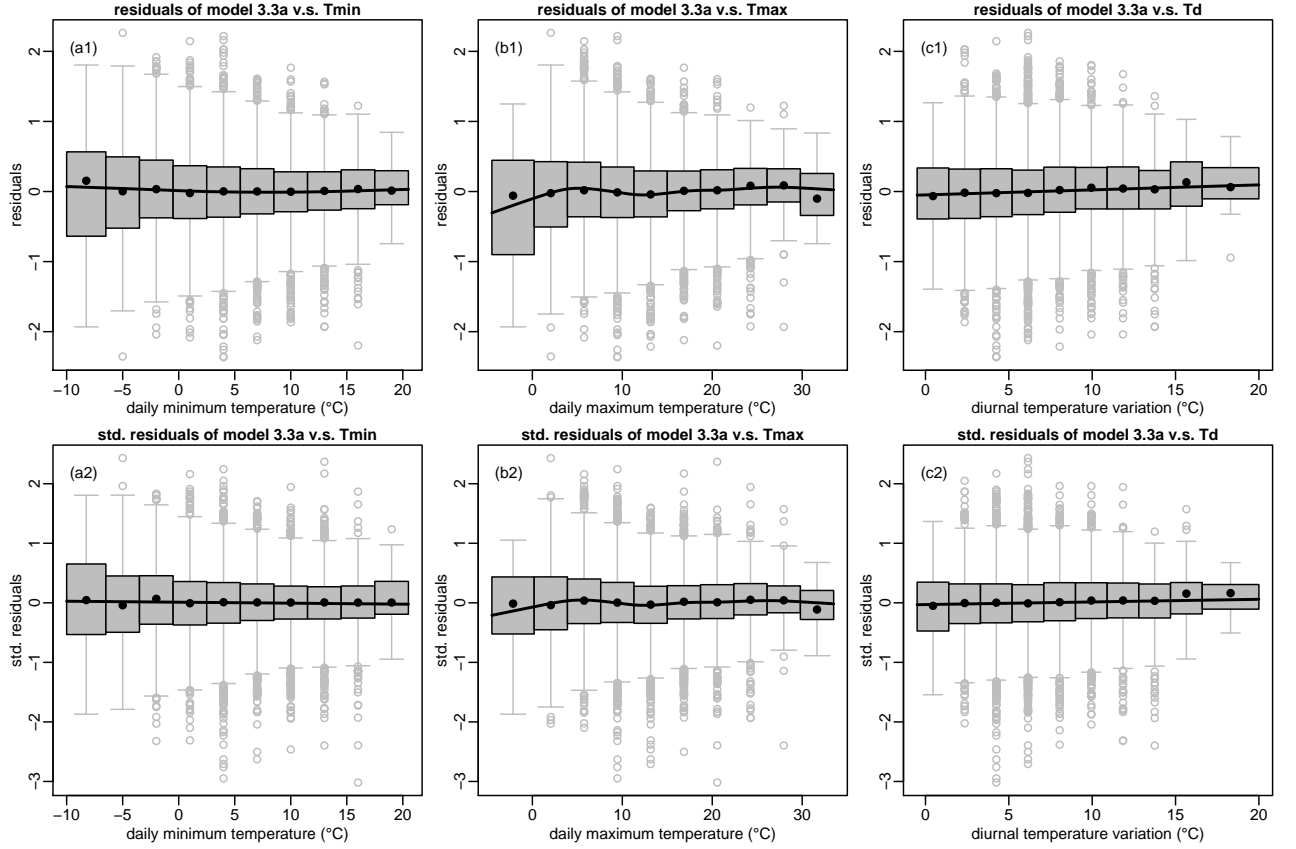
Fitting this model proves temperature variables extremely useful. They not only substantially reduces the estimated residual variance (model 3.2b has residual variance 0.3321138, while model 3.3a has residual variance 0.2799806), but also slightly decreases the autocorrelation in residuals (model 3.2b has estimated AR(1) correlation at 0.5447645, while the point estimate for model 3.3a is 0.517817). This is probably because that daily temperature is also a time series with strong autocorrelation, which can explain some autocorrelation in logBS as well. (Later in Figure 3.20 I will highlight how effectively temperature variables can improve model prediction.) Table 3.5 summarizes the effective degree of freedom for all functions in the model (note that  $f_5$  ends up with a value of 3.75 so it is slightly nonlinear!). Figure 3.10 once again inspects model residuals against temperature variables. It is interesting to notice that while daily maximum temperature has not been used for building model, the previously spotted unmodelled trend in panel (b1) and panel (b2) or Figure 3.9 has largely gone!

It is less easy to judge whether there is an interaction between  $\mathbf{T}_d^0$  and  $\mathbf{T}_d^*$ , that is, should the following model be considered,

**Model 3.3b:** model 3.3a +  $f_6$

$$\begin{aligned} \log \text{BS}_d = & u_1 + \sum_{s=2}^7 \delta_s^{\text{d}_w} u_s + f_1(\mathbf{w}_y; 15) + \sum_{s=2}^7 \delta_s^{\text{d}_w} f_{1,s}(\mathbf{w}_y; 15) + \\ & f_2(\mathbf{w}; 40) + \sum_{s=2}^7 \delta_s^{\text{d}_w} f_{2,s}(\mathbf{w}; 40) + f_3(\mathbf{w}_y, \mathbf{w}; 10, 40) + \\ & f_4(\mathbf{T}_d^0; 15) + f_5(\mathbf{T}_d^*; 10) + f_6(\mathbf{T}_d^0, \mathbf{T}_d^*; 15, 10) + e_d, \end{aligned}$$

where  $f_6$  is a tensor product spline between a natural cubic spline of  $\mathbf{T}_d^0$  and a natural cubic spline of  $\mathbf{T}_d^*$ ? Actually fitting this model shows that  $f_6$  has a degree of freedom at 20.11 (see Table 3.6)



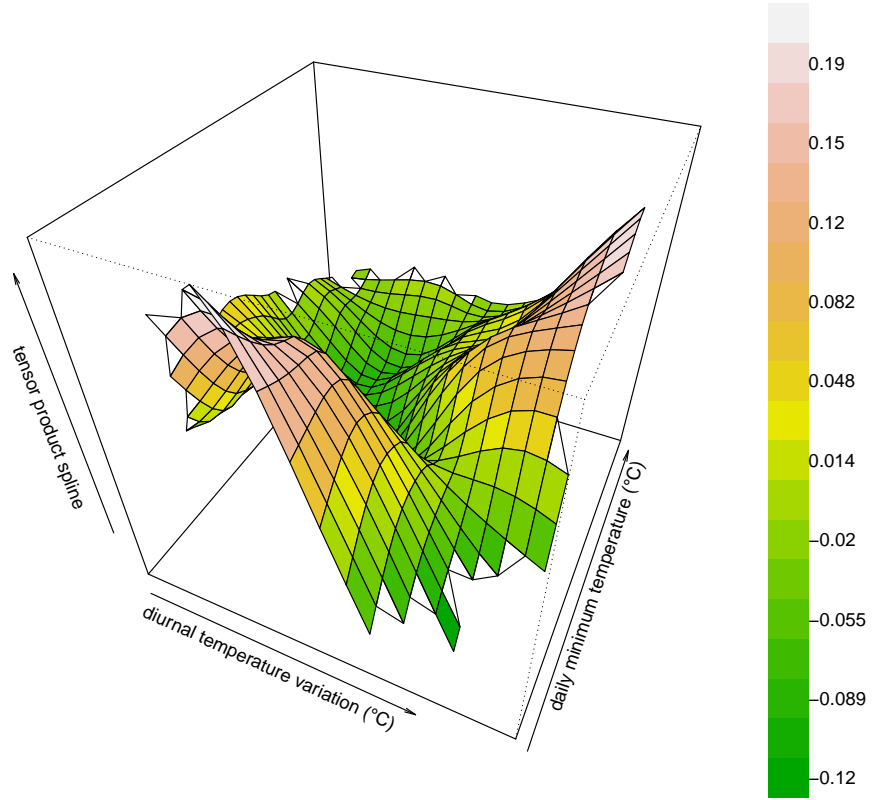
**Figure 3.10:** Inspecting residuals of fitted model 3.3a against daily minimum temperature (Tmin)  $T_d^0$ , daily maximum temperature (Tmax)  $T_d^1$  and diurnal temperature variation (Td)  $T_d^*$ . Unmodelled trend w.r.t.  $T_d^0$  and  $T_d^1$  which are previously spotted in Figure 3.9 for model 3.2b have been removed. The unmodelled trend w.r.t.  $T_d^*$  has also been largely taken off, although this variable is not used for building model 3.3a.

**Table 3.6:** Effective degree of freedom for all functions in fitted model 3.3b, assuming AR(1) error with unknown autocorrelation coefficient.

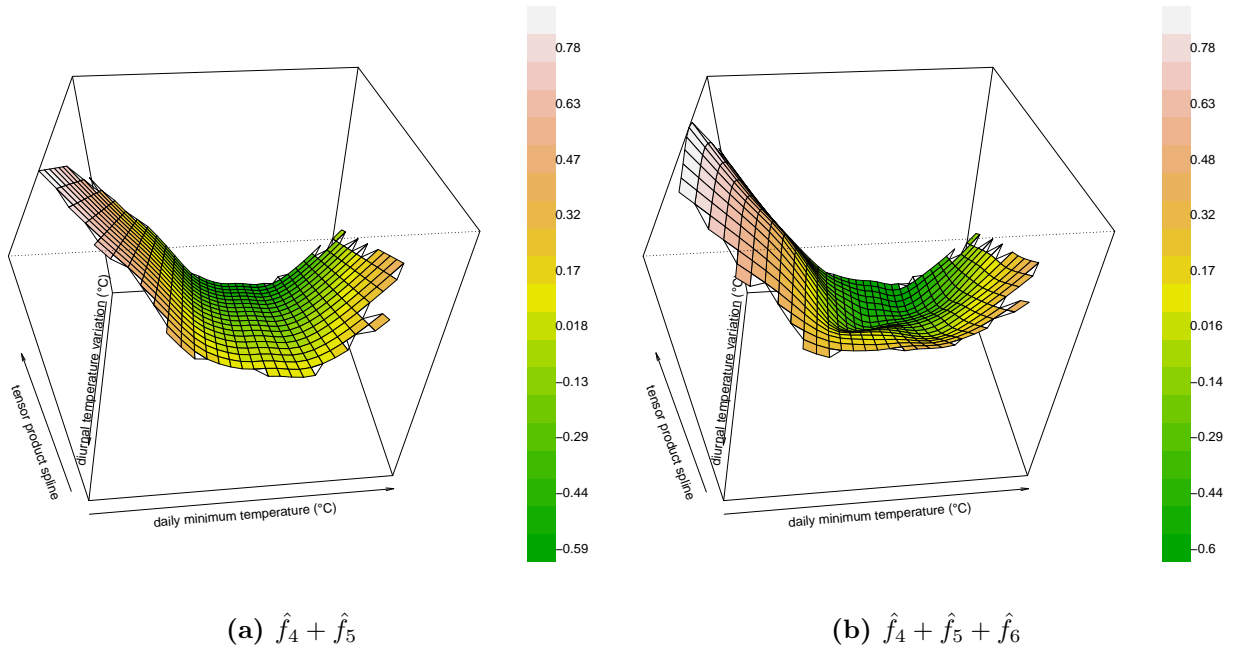
$\hat{f}_1$	$\hat{f}_{1,2}$	$\hat{f}_{1,3}$	$\hat{f}_{1,4}$	$\hat{f}_{1,5}$	$\hat{f}_{1,6}$	$\hat{f}_{1,7}$
9.69	0.00	0.00	0.00	1.10	1.33	2.17
$\hat{f}_2$	$\hat{f}_{2,2}$	$\hat{f}_{2,3}$	$\hat{f}_{2,4}$	$\hat{f}_{2,5}$	$\hat{f}_{2,6}$	$\hat{f}_{2,7}$
20.83	1.71	2.90	3.73	2.47	1.00	8.45
$\hat{f}_3$	$\hat{f}_4$	$\hat{f}_5$	$\hat{f}_6$			
86.85	7.67	1.00	20.11			

which is relatively high. Perhaps some idea can be obtained by surface visualization. Figure 3.11 is a perspective plot for  $\hat{f}_6$ ; it does appear that there is nonlinear interaction between  $T_d^0$  and  $T_d^*$ . Figure 3.12 compares surface of  $\hat{f}_4 + \hat{f}_5$  and surface of  $\hat{f}_4 + \hat{f}_5 + \hat{f}_6$ . The interaction does add some extra curvature. At this stage it is probably a good idea not to overly restrict the model structure, as eventually we need to build models for the whole monitoring network. So I will choose to retain this interaction term in the model for now.

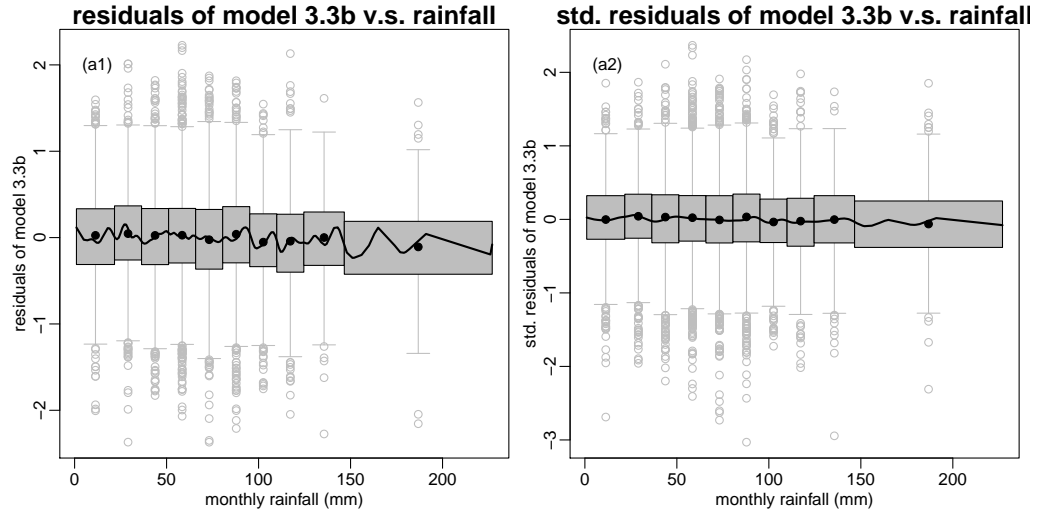
The UKCP09 dataset also provides rainfall data but only at a monthly resolution. It is unclear whether a monthly covariate is helpful in explaining daily logBS. Figure 3.13 sketches residuals and standardized residuals of fitted model 3.3b against monthly rainfall. It appears that at least at station *Manchester 11*, there isn't strong relationship between rainfall and logBS. However, this can not exclude the impact of rainfall in general. Notably, the rainfall measurements at this particular station only vary between 1 and 227, however, rainfall measurements at all stations of the Network vary between 1 and 583 which is a much wider range. So the effect of rainfall needs be investigated at a later stage.



**Figure 3.11:** Perspective plot for  $\hat{f}_6$  in fitted model 3.2b. The surface is highly nonlinear, so there is probably a notable interaction between  $T_d^0$  and  $T_d^*$ . See §2.3 for more details on such kind of plot.



**Figure 3.12:** Comparing  $\hat{f}_4 + \hat{f}_5$  in the fitted model 3.3a and  $\hat{f}_4 + \hat{f}_5 + \hat{f}_6$  in the fitted model 3.3b. The interaction adds some extra, noticeable curvature to the surface.



**Figure 3.13:** Inspecting residuals and standardized residuals of fitted model 3.3b against monthly rainfall. It appears that at station *Manchester 11*, there isn't strong relationship between rainfall and logBS. However, this can not exclude the impact of rainfall in general. Notably, the rainfall measurements at this particular station only vary between 1 and 227, however, rainfall measurements at all stations of the Network vary between 1 and 583 which is a much wider range. So the effect of rainfall needs to be investigated at a later stage.

### 3.2.5 Model checking and visualization

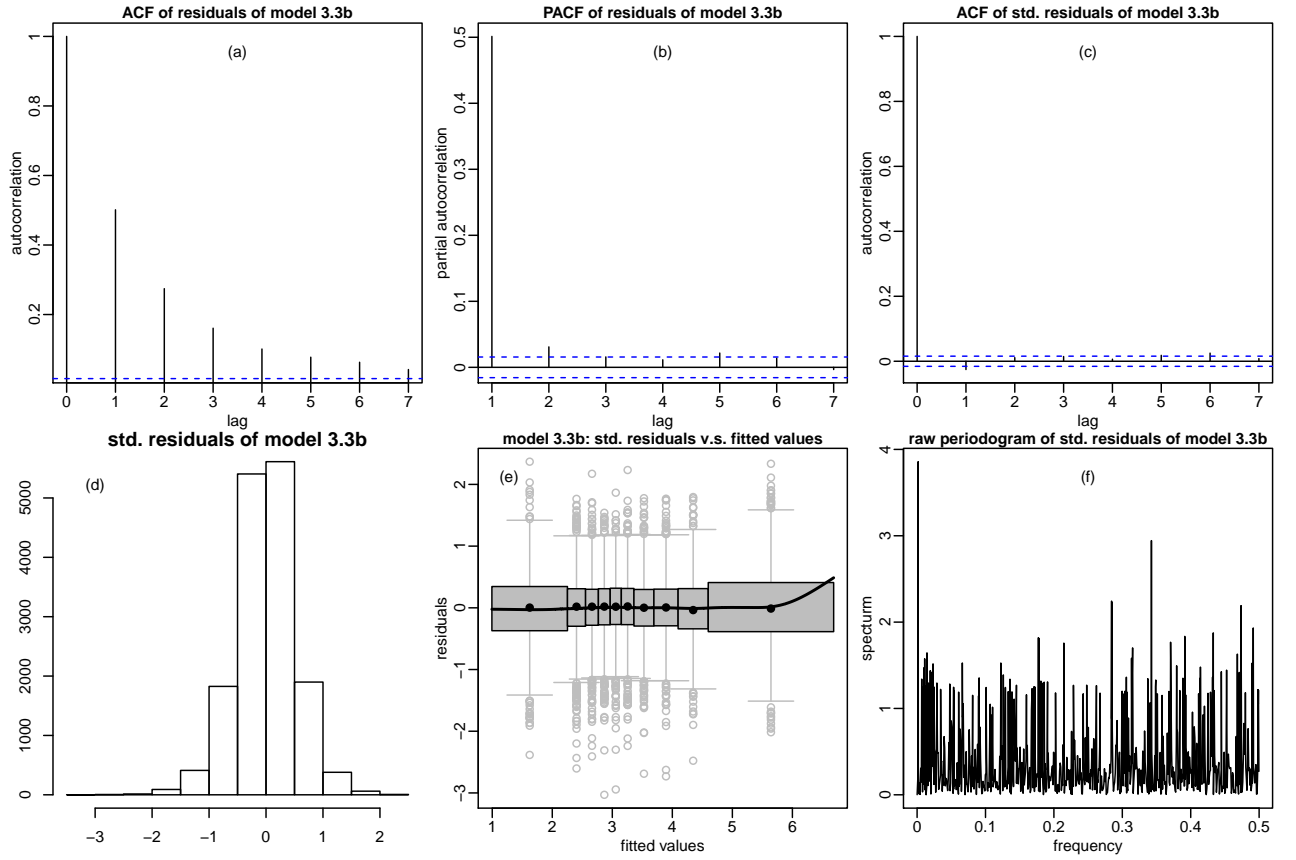
The fitted model 3.3b involves 732 coefficients for 15707 data, with a resulting degree of freedom at 208.64. The estimated residual variance is 0.2779, while the sample variance of logBS observations is 1.1144, so the fitted model has an adjusted R-squared of  $1 - 0.2779 / 1.1144 = 75.06\%$ . I will now perform some model checking for this model and visualize all spline functions.

Figure 3.14 conducts some residual checking. Panels (a) to (c) just reassure that autocorrelation in daily logBS has been adequately modelled by an AR(1) process. Panel (e) for residuals v.s. fitted values might cause some concern. However, the increase at the right boundary really implies that a smoothing model has difficulty in approximating unextraordinarily large observations (extreme values). The highest value observed at *Manchester 11* is 8.6 for logBS, but  $5413 \mu\text{gm}^{-3}$  in the original scale (even larger than the  $4500 \mu\text{gm}^{-3}$  in London's 1952 episode of fog!). The fitted model predicts 6.36 for logBS and the residual is as high as 2.23. In general, the constant error variance assumption is not violated.

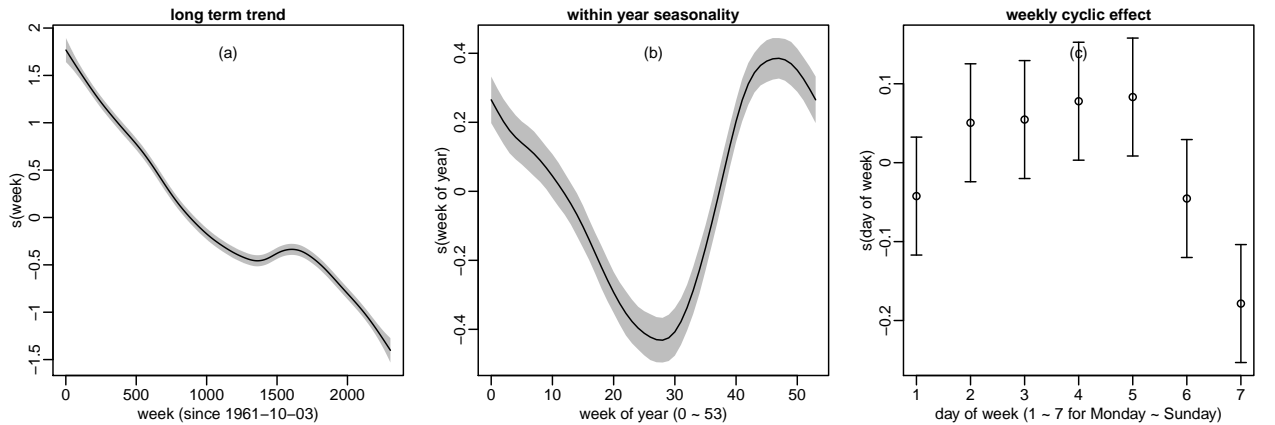
Figure 3.15 is an illustration of the estimated long term trend  $\hat{f}_2$ , within-year seasonality  $\hat{f}_1$  and weekly effect  $\hat{u}_1, \hat{u}_1 + \hat{u}_2, \dots, \hat{u}_1 + \hat{u}_7$ . There are a few interesting observations.

1. The long term trend follows a general decline, except for some time around week 1500 (around mid 1980's).
2. The "V" shape curve of seasonal pattern is not surprising, as winter is when heating demands and fuel combustion are high. However, the peak of the curve is somewhere around November not December. Perhaps people went on holiday near Christmas and there is less industrial and commercial activity than usual?
3. The weekly pattern suggests that Sunday is normally when Smoke concentration is the lowest in a week. However, the difference between Saturday and Sunday is rather notable. Also, while Tuesday, Wednesday, Thursday and Friday are more alike, Monday is quite different. The reason behind this phenomenon is unknown, but it is readily clear that it is not a good idea to simply model weekly effect as a two-level factor for weekday / weekend.

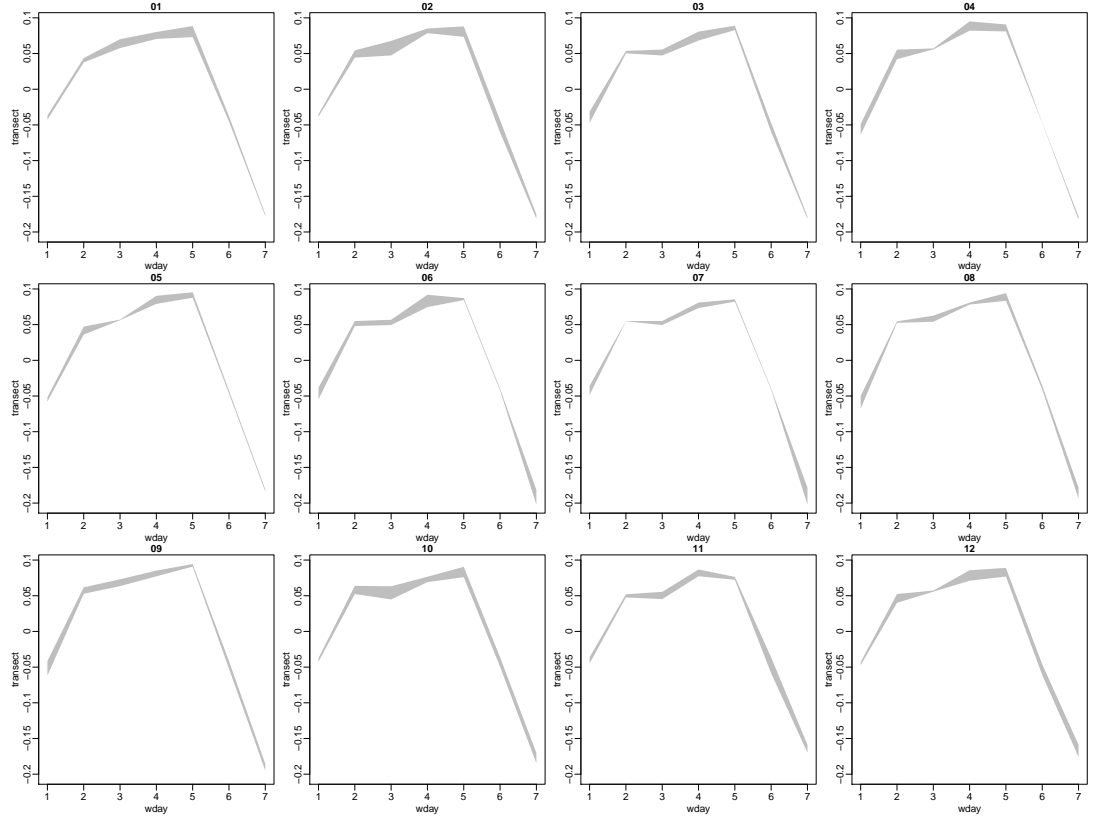
The most obvious way to visualize interaction term  $\hat{f}_{1,2}, \dots, \hat{f}_{1,3}$  is to plot them individually against



**Figure 3.14:** Residual checking for fitted model 3.3b. Panel (a) is the ACF for residuals. Panel (b) is the PACF for residuals. Panel (c) is the ACF for AR(1) standardized residuals. These panels convince that autocorrelation in daily logBS has been well modelled. Panel (d) is a histogram of standardized residuals. Panel (e) sketches standardized residuals against fitted values. Generally the constant variance assumption for model error is not violated. The increase at the right boundary implies that a smoothing model has difficulty in approximating unextraordinarily large observations (extreme values). The highest value observed at *Manchester 11* is 8.6 for logBS, but  $5413 \mu\text{gm}^{-3}$  in the original scale. The fitted model predicts 6.36 for logBS and the residual is as high as 2.23. Panel (f) is a raw periodogram for standardized residuals. It verifies that they are almost white noise without any unmodelled cyclic / periodic effect.



**Figure 3.15:** Illustration of the general long term trend  $\hat{f}_2$ , within-year seasonality  $\hat{f}_1$  and weekly effect  $\hat{u}_1, \hat{u}_1 + \hat{u}_2, \dots, \hat{u}_1 + \hat{u}_7$  in fitted model 3.3b. The gray ribbon covers region that is within 2 standard deviation of the estimate.



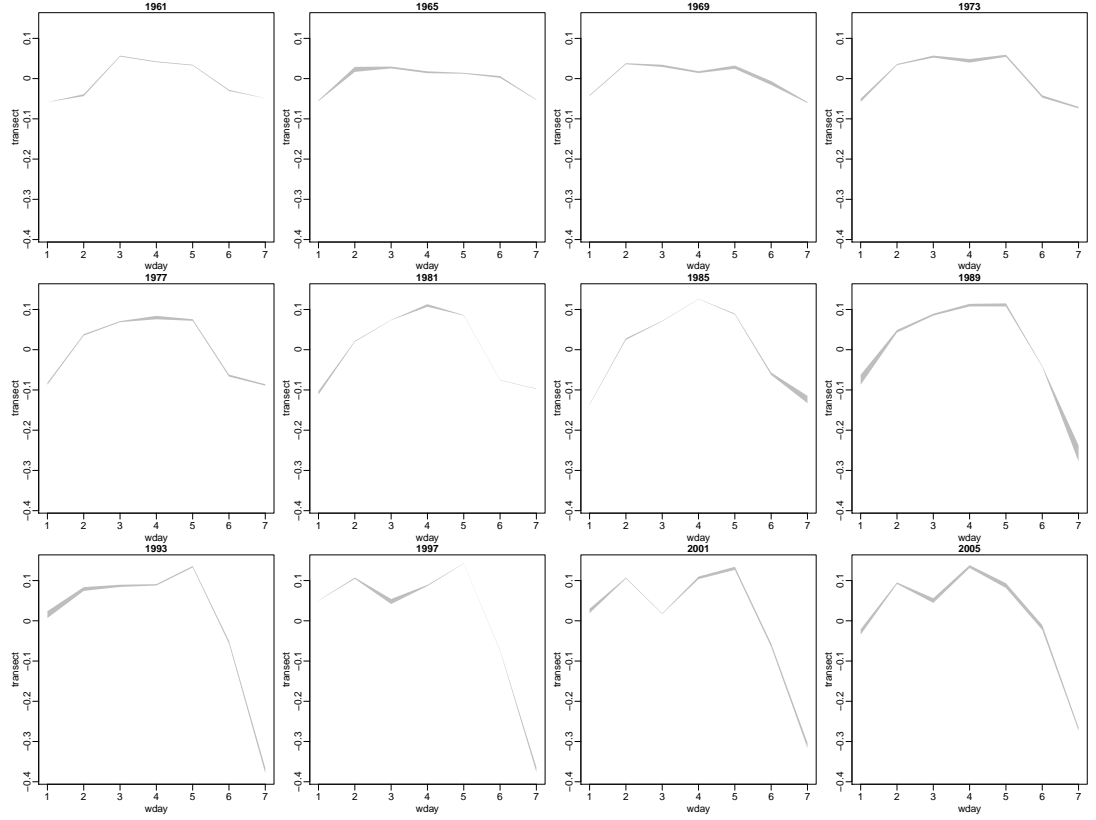
**Figure 3.16:** Grouped transect plots (by month) for  $\hat{f}_1$ ,  $\hat{f}_1 + \hat{f}_{1,2}$ , ...,  $\hat{f}_1 + \hat{f}_{1,7}$  in the fitted model 3.3b. There is plenty of **cyclic** variation of weekly pattern with a year, notably for the middle of a week (i.e., Tuesday to Friday).

$\mathbf{w}_y$ . Similarly  $\hat{f}_{2,2}$ , ...,  $\hat{f}_{2,7}$  can be each plotted against  $\mathbf{w}$ . However, it is probably more informative to learn how the weekly pattern varies within a year and changes in a long term. For this purpose, I will for example, fix  $\mathbf{w}_y$  at a specific value and plot the resulting 7 values  $\hat{f}_1(\mathbf{w}_y)$ ,  $\hat{f}_1(\mathbf{w}_y) + \hat{f}_{1,2}(\mathbf{w}_y)$ , ...,  $\hat{f}_1(\mathbf{w}_y) + \hat{f}_{1,7}(\mathbf{w}_y)$ . Yet this is not ideal as this single line does not express the local variation along the other coordinate  $\mathbf{w}_y$ . A better option is to group transects and plot them in chunks. For example, for  $\hat{f}_1(\mathbf{w}_y)$ ,  $\hat{f}_1(\mathbf{w}_y) + \hat{f}_{1,2}(\mathbf{w}_y)$ , ...,  $\hat{f}_1(\mathbf{w}_y) + \hat{f}_{1,7}(\mathbf{w}_y)$  I will group weeks of year into 12 months and produce 12 grouped transect plots. For  $\hat{f}_2(\mathbf{w})$ ,  $\hat{f}_2(\mathbf{w}) + \hat{f}_{2,2}(\mathbf{w})$ , ...,  $\hat{f}_2(\mathbf{w}) + \hat{f}_{2,7}(\mathbf{w})$  I will group weeks into 45 years and produce 45 grouped transect plots. For  $\hat{f}_3(\mathbf{w}_y, \mathbf{w})$ , the interest is how seasonal pattern would vary in a long term and the idea of grouped transect plots applies as well. Figure 3.16, Figure 3.17 and Figure 3.18 illustrate such grouped transect plots for these model components.

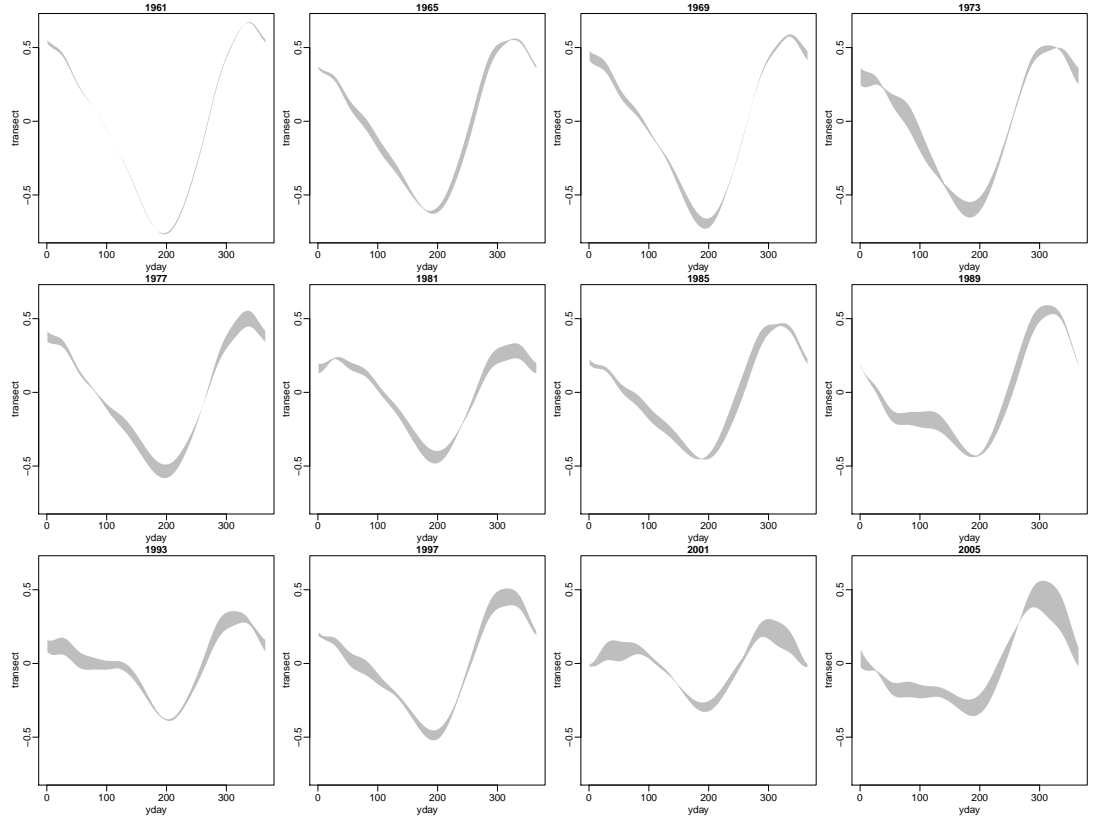
1. From Figure 3.16 we see that there is plenty of **cyclic** variation of weekly pattern within a year (notably for the middle of a week, i.e., Tuesday to Friday), that is, these changes will be replicated for all 45 years. .
2. From Figure 3.17 we see that there isn't much **non-cyclic** variation of weekly pattern within a year (the gray ribbon contains 53 transect lines, but they are pretty much clustered together; so this can probably be interpreted that all variation of weekly pattern in a year is cyclic, as is captured by  $\hat{f}_{1,2}(\mathbf{w}_y, \mathbf{w})$  to  $\hat{f}_{1,7}(\mathbf{w}_y, \mathbf{w})$ ), but the long term change is dramatic.
3. From Figure 3.18 we see that the seasonal patter can vary greatly within a year. This needs some clarification. Isn't a year a period of seasonal pattern? How can it vary? Well, the general seasonality  $\hat{f}_1$  is static, but it does not mean that it is simply replicated for all 45 years, or only rescaled once per year; its shape can change smoothly as each day passes. So this tensor product spline is capable to model very complicated dynamics.

In the previous section we have visualized the interaction term  $\hat{f}_6$  for  $\mathbf{T}_d^0$  and  $\mathbf{T}_d^*$ . Their marginal effects are now illustrated in Figure 3.19. There should be some more sophisticated scientific explanation on

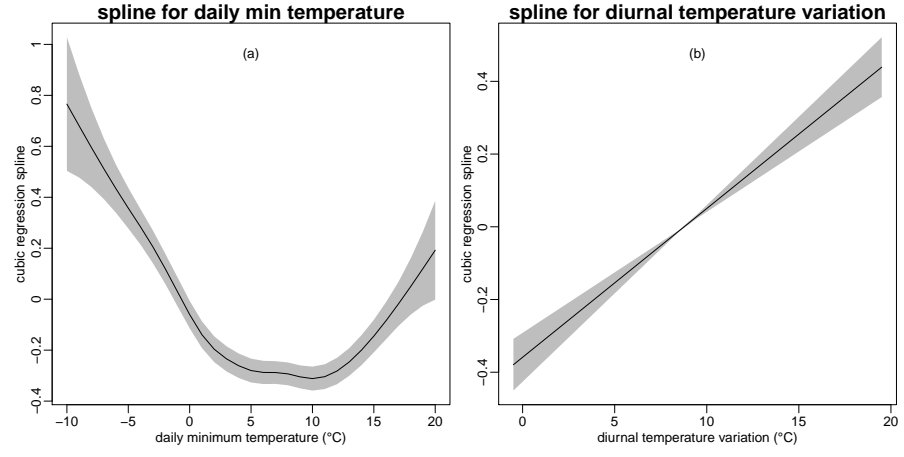




**Figure 3.17:** Grouped transect plots (by year) for  $\hat{f}_2$ ,  $\hat{f}_2 + \hat{f}_{2,2}$ , ...,  $\hat{f}_2 + \hat{f}_{2,7}$  in the fitted model 3.3b. There isn't much **non-cyclic** variation of weekly pattern within a year (the gray ribbon contains 53 transect lines, but they are pretty much clustered together), but the long term change is dramatic. Note that for space reason, the plot is only made every 4 years.



**Figure 3.18:** Grouped transect plots (by year) for  $\hat{f}_1(\mathbf{w}_y) + \hat{f}_3(\mathbf{w}_y, \mathbf{w})$  in the fitted model 3.3b. The wide width of the gray ribbon implies that the seasonal pattern vary greatly within a year. For space reason, the plot is only presented every 4 years.



**Figure 3.19:** Marginal effects of daily minimum temperature ( $\hat{f}_4$ ) and diurnal temperature variation ( $\hat{f}_5$ ) in fitted model 3.3b. There should be some more sophisticated (and interesting) scientific explanation on the impact of temperature on Smoke concentration, but it is out of scope of this PhD thesis.

the impact of temperature on Smoke concentration, other than the association between temperature and heating. For example, the positive slope of  $T_d^*$  is not something that can be explained by heating. However, this scientific topic is out of scope of this PhD thesis.

After visualizing each individual splines, let us see how they help prediction when added together. Figure 3.20 breaks down components of fitted model 3.3b into trend, seasonal, weekly and temperature components and illustrates their contribution to logBS prediction in four decades. The weekly component have greater and greater effects in later decades, and the temperature components substantially improve the quality of prediction.

### 3.2.6 Summary

The model development for a logBS time series has been very much detailed. Before moving on, I will summarize key messages that can be learnt from this *Manchester 11* case study.

First of all,

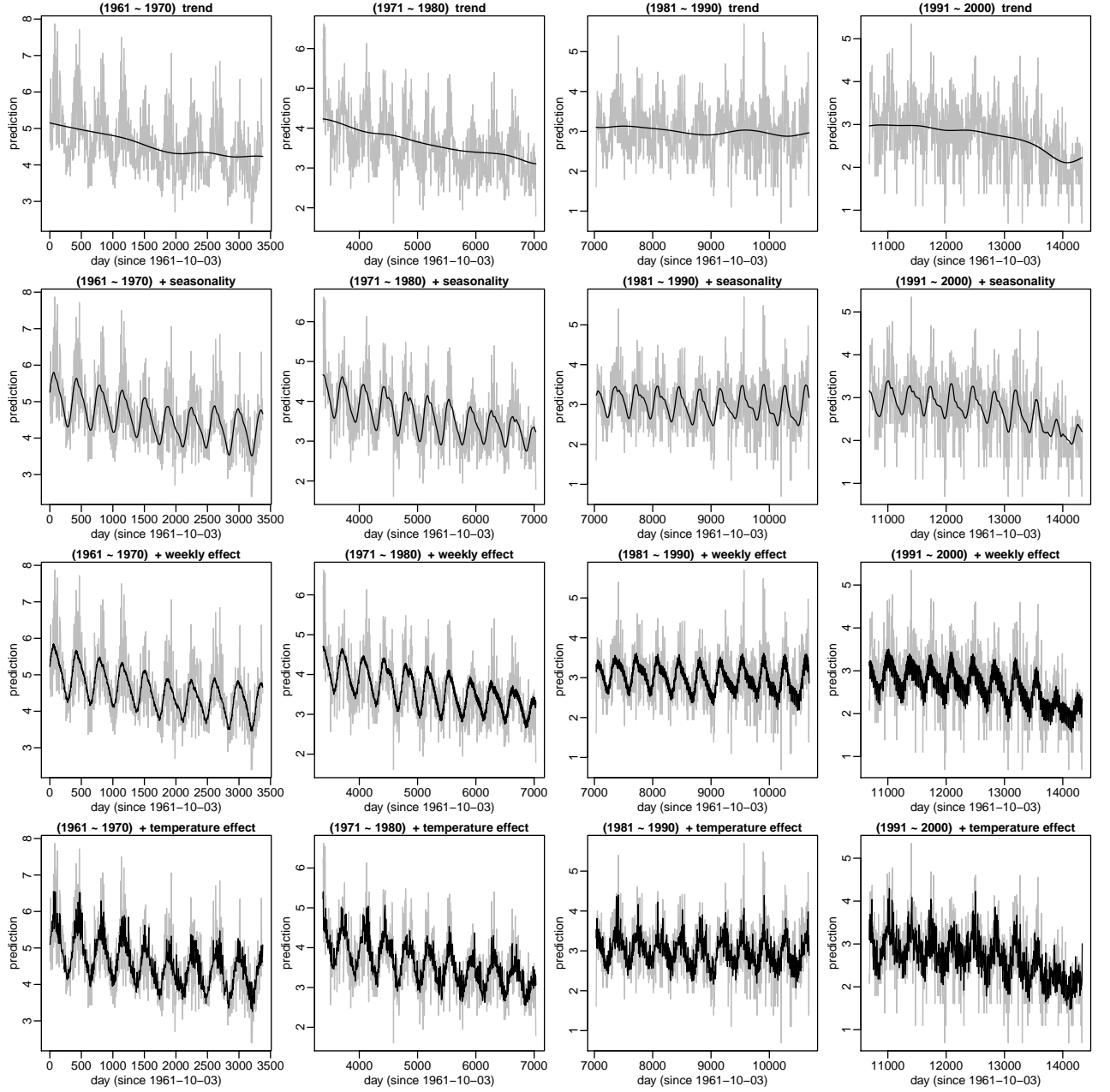
- A logBS daily time series can be decomposed into additive smooth functions of time variables at different resolutions (as well as their interactions), namely week  $w$ , day of week  $d_w$  and week of year  $w_y$ , plus an AR(1) process mainly for autocorrelation with a week.

In addition, we have obtained some rough idea of the complexity for each component:

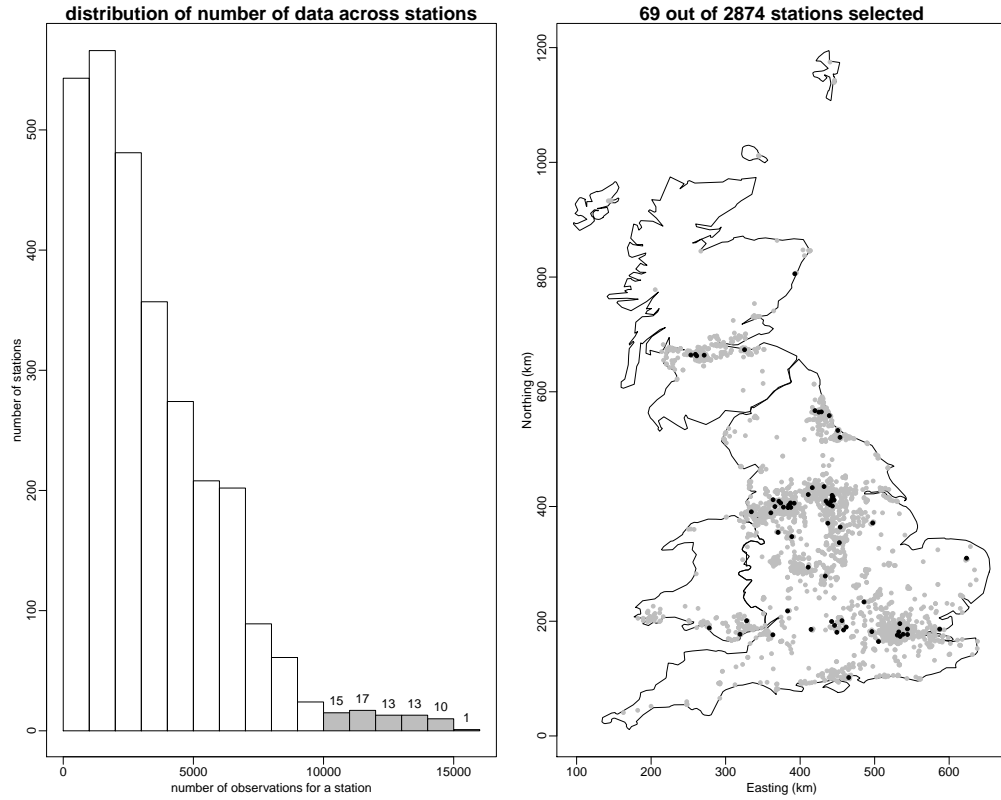
- Conditional on an adequate AR(1) model for autocorrelation in daily data, 15 ~ 20 knots should be sufficient for variable  $w_y$ , and 40 ~ 50 knots should be more than sufficient for variable  $w$ . Temperature variables do not need very high  $k$  values; 15 knots should be more than sufficient for both  $T_d^0$  and  $T_d^*$ .

Of course, it would be good if such rough idea drawn from *Manchester 11* case study can be validated on other stations.

As one way to do this, let us consider a subset of the full logBS dataset, where only stations with long monitoring history, say over 10000 days' records, are retained. It turns out that 69 stations are



**Figure 3.20:** An illustration of how the fitted model 3.3b predicts logBS when its model terms are added together. The daily time series is broken into four decades, each per column panel. For column panels: the topmost row shows trend component  $\hat{u}_1 + \hat{f}_2$  only; then the second row adds seasonal components  $\hat{f}_1$  and  $\hat{f}_3$ ; the third row further adds weekly components  $\hat{u}_2, \dots, \hat{u}_7, \hat{f}_{1,2}, \dots, \hat{f}_{1,7}$  and  $\hat{f}_{2,1}, \dots, \hat{f}_{2,7}$ ; finally the bottom row adds temperature components  $\hat{f}_4, \hat{f}_5$  and  $\hat{f}_6$ . The weekly components have increasingly noticeable effects as time goes by, and the temperature components substantially improve prediction.



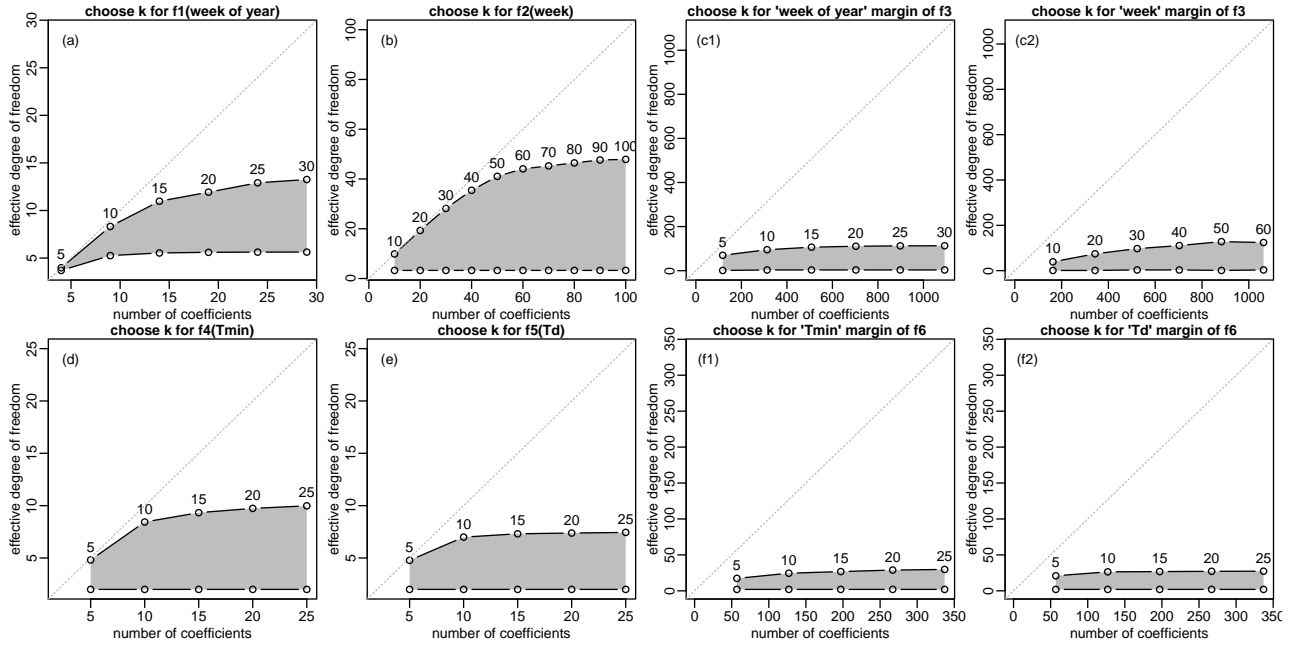
**Figure 3.21:** All 69 stations with over 10000 days’ measurements are selected to validate my rough estimate on the number of knots for various splines in model 3.3b. The panel on the left is a histogram for the number of measurements per station. Only a small proportion of stations have a long monitoring history (10000 days  $\approx$  30 years). The right panel identifies those stations on the map. From their relative positions (black dots) to all other historical stations (gray dots) in the Network, they seem to be pretty representative of the Network. These stations contribute to about 9% of the all daily measurements from the Network.

selected (see Figure 3.21). From their relative positions to all other stations in the Network, they seem to be pretty representative of the Network. Dropping weekly components from model 3.3b, we can fit the reduced model for logBS on each day of week per station (assuming i.i.d. errors of course). This ends up with  $69 \times 7 = 483$  models. Repeating the “choose k” procedure for all of them, I summarize the result in Figure 3.22. The Figure is produced in a similar way to Figure 3.4, except that the line plot is replaced by their bounding polygon. The effective degree of freedom plateaus for all splines as  $k$  increases, and except that I might have slightly underestimated the complexity required for  $f_2(\mathbf{w})$ , the estimates for all terms are fairly adequate.

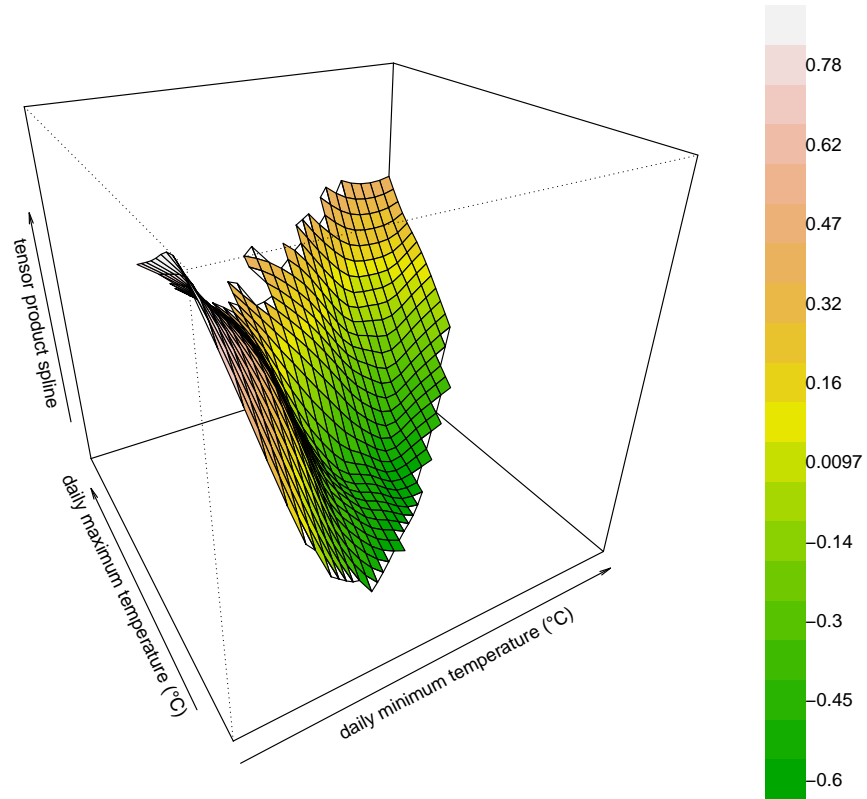
Another useful message is that meteorological variables, particularly temperature variables  $T_d^0$  and  $T_d^*$ , are extremely helpful for prediction. However, from the *Manchester 11* dataset (too small a sample size) the following can not be learned so they need be revisited in future.

1. Is the temperature effect variable in a long term, that is, is there an interaction between temperature and year  $y$ ?
2. Will rainfall affect logBS?

There are other alternative ways to model temperature effects. Using a bivariate thin-plate spline for  $T_d^0$  and  $T_d^*$  turns out to produce the same surface as is graphed in Figure 3.12. In fact, it is also equivalent to use a bivariate thin-plate spline for  $T_d^0$  and  $T_d^1$ , although the surface lies in a different domain (see Figure 3.23). The thin-plate spline representation is very concise: we can just write one function instead of three for temperature variables. However, it is probably preferable to use the tensor product spline representation:  $f_4(T_d^0) + f_5(T_d^*) + f_6(T_d^0, T_d^*)$  where the marginal effects  $f_4$  and  $f_5$  can be visualized (see Figure 3.19).



**Figure 3.22:** The result of “choose  $k$ ” process for fitting model 3.3b (replacing  $d$  with  $v$ , dropping weekly components and hence assuming i.i.d. errors) for logBS on each day of week per station. This gives a good estimate on adequate  $k$  values for the model.



**Figure 3.23:** Using a bivariate thin-plate spline for  $T_d^0$  and  $T_d^*$  produces a surface on a different domain. Since these two variables are correlated the domain stretches along the diagonal of the square. But after a transformation  $T_d^0 = T_d^1$  and  $T_d^* = T_d^1 - T_d^0$  the surface almost coincides the surface in Figure 3.12.

There are also other ways to introduce  $\mathbf{d}_w$  into the model than using a factor and replicating the smooth function for each level. These alternatives include:

- using a cubic cyclic spline for  $\mathbf{d}_w$ ;
- using a random effect (penalized factor) for  $\mathbf{d}_w$ .

In both cases the interaction between  $\mathbf{d}_w$  and other variables can be constructed as a tensor product spline. For example, model 3.2b may be written as

$$\begin{aligned} \log\text{BS}_d = & f_0(\mathbf{d}_w; 7) + f_1(\mathbf{w}_y; 15) + \tilde{f}_1(\mathbf{d}_w, \mathbf{w}_y; 7, 15) + \\ & f_2(\mathbf{w}; 40) + \tilde{f}_2(\mathbf{d}_w, \mathbf{w}; 7, 40) + f_3(\mathbf{w}_y, \mathbf{w}; 10, 40), \end{aligned}$$

which is a great simplification as all “ $\sum$ ” are gone. From now on, this tensor product spline construction will be used for modelling interaction with  $\mathbf{d}_w$  (and generally any interaction between a spline and a factor).

### 3.3 Motivating spatial-temporal modelling

In the previous section I have demonstrated how to build a time series model for daily logBS from any station, and in fact, as I have verified, for all 69 representative stations marked in Figure 3.21, the proposed model 3.3b proves adequate.

However, no matter how many individual daily time series are fitted and how good the prediction is at those stations, it is not possible to predict logBS away from those stations. To make such spatial or spatial-temporal prediction, we need to assume that those logBS times series are correlated, or that they vary smoothly over space hence a surface prediction is possible.

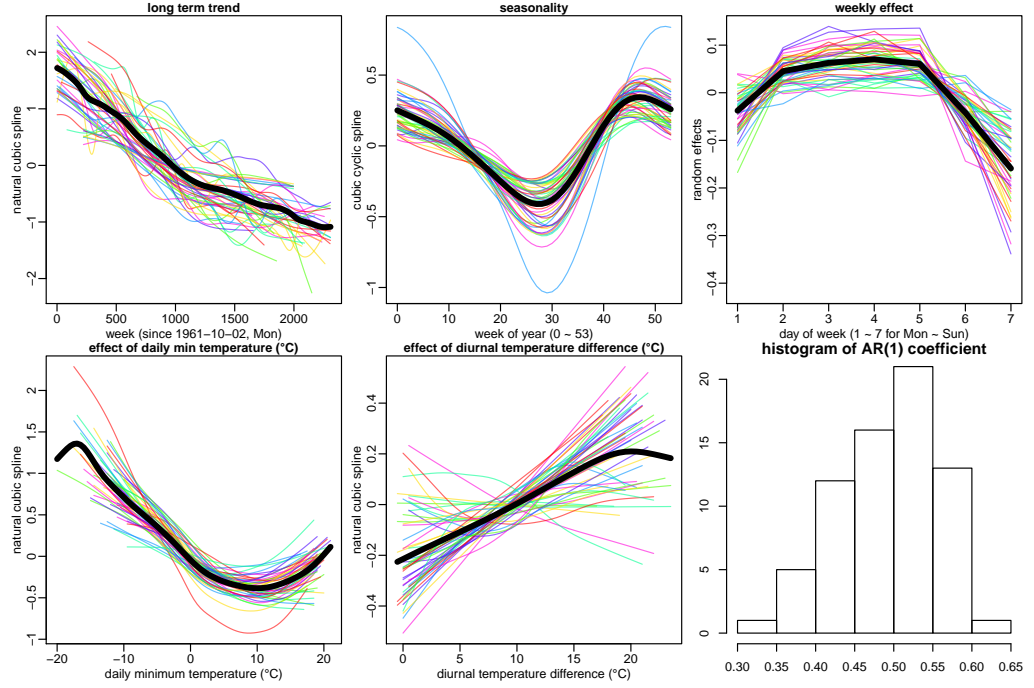
Figure 3.24 illustrated in coloured curves the estimated main effects  $\hat{f}_0$ ,  $\hat{f}_1$ ,  $\hat{f}_2$ ,  $\hat{f}_4$  and  $\hat{f}_5$  in fitted model 3.3b for all 69 stations. It can be observed that they vary greatly between stations. However, the general pattern of each effect seems very simple, as is highlighted by the bold black curve on each graph. In spatial-temporal modelling, it is assumed that each individual time series is an observation from a spatial population. While the spatial variability between observations can be high and complicated, the mean of the population is smooth and simple.

Additive models via penalized regression splines can be used for such spatial-temporal modelling. The population mean can be modelled by a smooth function, like a cubic spline of week variable  $\mathbf{v}$  for long term trend and a thin-plate spline of spatial coordinates  $\{\mathbf{e}_i, \mathbf{n}_i\}$  for spatial trend. The spatial variability across time can be modelled by a tensor product spline with a thin-plate spline margin as its spatial component and a cubic spline margin as its temporal component. In the next section, I will demonstrate this modelling approach for annual mean logBS from all stations.

## 3.4 Spatial-temporal modelling for yearly mean logBS

### 3.4.1 Description of data

The spatial-temporal dataset for yearly mean logBS, is obtained from aggregating daily logBS data per year and station. However, to ensure that this sample mean is representative for a whole year,



**Figure 3.24:** Estimated main effects (thin coloured curves) in fitted model 3.3b for all 69 stations (described in Figure 3.21) and their mean (in a bold black curve). Spatial-temporal modelling assumes each colour curve to be a sample from a population with the bold black curve as its mean. The variability between a sample and the mean is to be explained spatially: samples observed at closer spatial locations are more similar than those observed at far apart locations. The final histogram is for the estimated AR(1) coefficients in each individual time series model.

aggregation is not done if there are fewer than 273 (approximately 75% of number of days in a year) daily data. The resulting dataset has 24239 data from a total number of 2626 stations, covering a period from year 1962 to 2004. See Figure 3.25 and Figure 3.26 for some illustration of this dataset.

### 3.4.2 A basic model

Shaddick and Zidek (2014) attempted spatial-temporal modelling of such dataset on a study period 1966 ~ 1996, concluding that there is a strong interaction between space and time (i.e., the spatial variability of logBS varies in time)<sup>1</sup>. A linear mixed model with the following core structure was suggested:

$$\log BS_{iy} = \beta_0 + \beta_1 y + \beta_2 y^2 + \beta_{0i} + \beta_{1i} y + \beta_{2i} y^2 + \epsilon_{iy}.$$

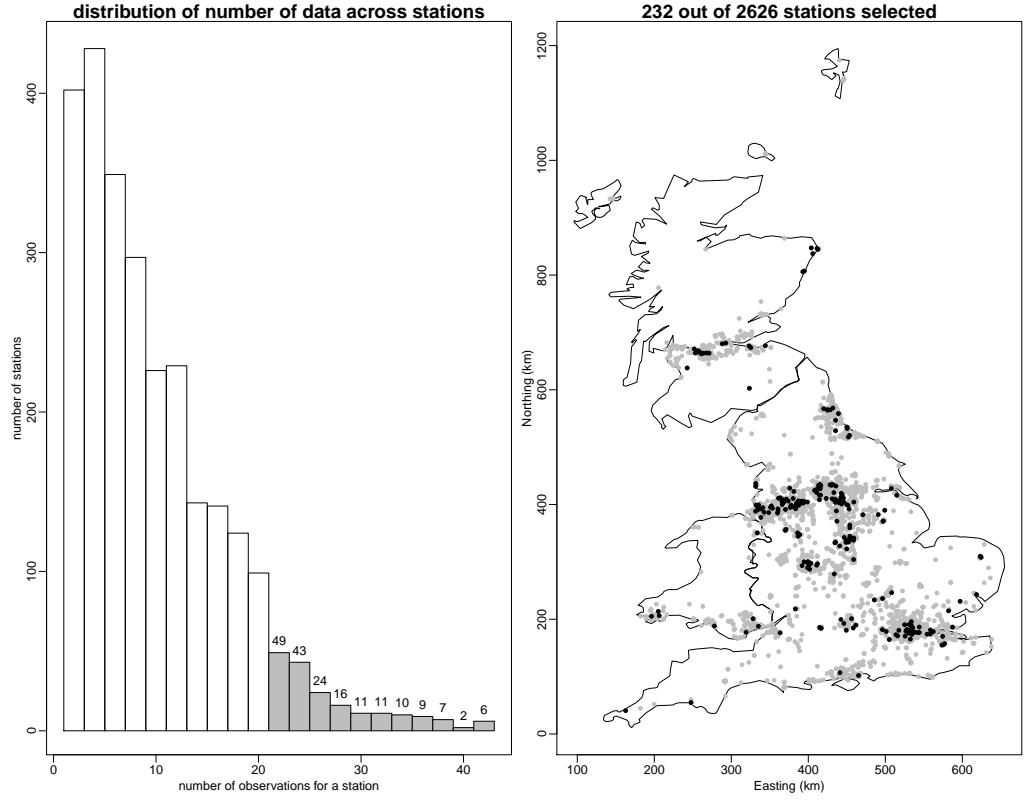
The long term trend is modelled by a quadratic polynomial of  $y$ , whose coefficients are allowed to vary spatially. Specifically,  $\beta_0$ ,  $\beta_1$  and  $\beta_2$  are fixed effect for the mean polynomial coefficients over space, and  $\beta_{0i}$ ,  $\beta_{1i}$  and  $\beta_{2i}$  are zero-mean spatial random effects with independent Matérn covariance. Representing this model using an additive model with splines is a type of polynomial “by” model (see model 2.2) which I have previously introduced in 2.1.1:

$$\log BS_{iy} = f_1(\{\mathbf{e}_i, \mathbf{n}_i\}) + f_2(\{\mathbf{e}_i, \mathbf{n}_i\})y + f_3(\{\mathbf{e}_i, \mathbf{n}_i\})y^2 + \epsilon_{iy},$$

where  $f_1$ ,  $f_2$  and  $f_3$  are thin-plate splines of  $(\{\mathbf{e}_i, \mathbf{n}_i\})$ , modelling spatially vary polynomial coefficients.

A quadratic polynomial is unlikely to adequately model trend at a station. For example, the trend at *Manchester 11* shown in Figure 3.15 is much more complicated than a polynomial. Furthermore, in §2.1.1 it has been explained that a polynomial “by” model suffers from Runge’s phenomenon. A

<sup>1</sup>In fact, the long term trend plot in Figure 3.24 is also an evidence for this.



**Figure 3.25:** [Illustration of annual mean logBS dataset (part 1)] As we have learnt from the histogram in Figure 3.21, the majority of stations in the Network did not operate for long enough time compared with the Network itself. The histogram on the left of this Figure tells the same thing. Among 2626 stations (gray dots in the map on the right panel) in the annual logBS dataset, only 232 of them (black dots on the map) had operated for over 20 years. They contributed to 25% of all annual data.

both more flexible and numerically stable model representation is to use tensor product splines. Let us consider the following model:

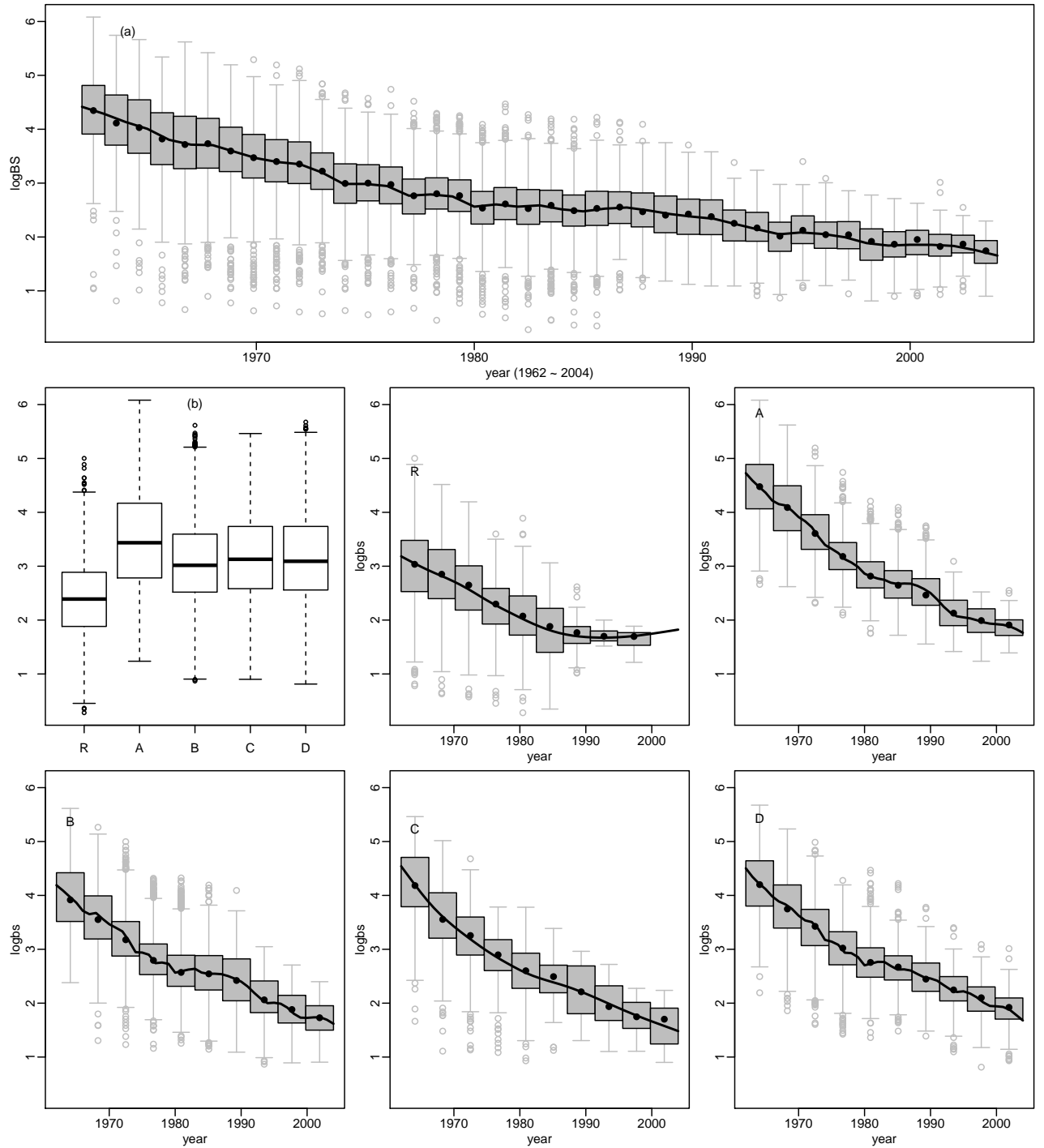
**Model 3.4:**

$$\log\text{BS}_{iy} = f_0(\mathbf{E}_i; 5) + f_1(y; k_1) + \tilde{f}_1(\mathbf{E}_i, y; 5, k_1) + f_2(\{\mathbf{e}_i, \mathbf{n}_i\}; k_2) + f_3(y, \{\mathbf{e}_i, \mathbf{n}_i\}; k_{31}, k_{32}) + f_4(\mathbf{h}_i; k_4) + \epsilon_{iy},$$

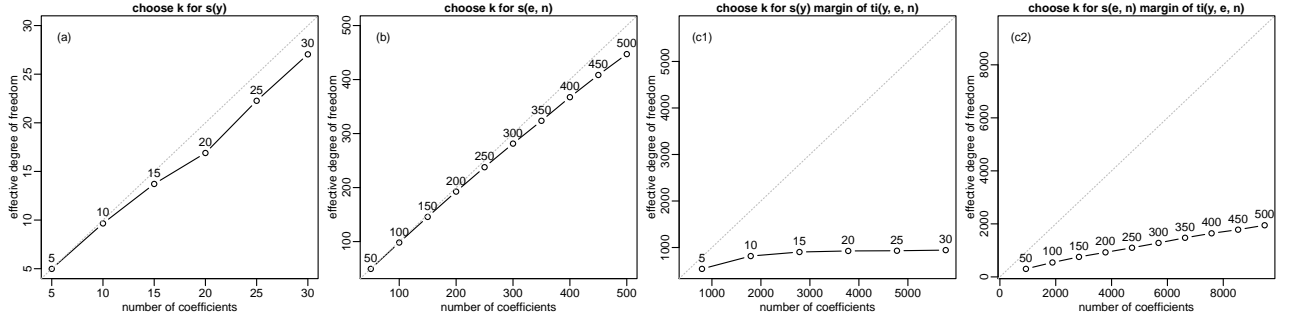
where

- $f_0$  is a random effect modelling the mean of logBS from stations of environmental type  $\mathbf{E}_i$ ;
- the general long term trend (common to all stations) is modelled by a natural cubic spline  $f_1$  of  $y$  with  $k_1$  knots;
- $\tilde{f}_1$  is a tensor product spline, whose first margin is a random effect for  $\mathbf{E}_i$  and the second margin is cubic spline of  $y$  with  $k_1$  knots. This component models how  $f_1$  varies for each  $\mathbf{E}_i$ ;
- the general spatial distribution of logBS across stations (invariant with time) is modelled by a thin-plate spline  $f_2$  of  $\{\mathbf{e}_i, \mathbf{n}_i\}$  with rank  $k_2$  (thin-plate spline is a knots-free spline; it places knots exact all all unique locations; the  $k$  value associated with a thin-plate spline is the rank used in its low rank approximation; see §1.1.2 if you need a revision);
- the time varying spatial variability is modelled by a tensor product spline  $f_3$ , in which the first margin is a natural cubic spline of  $y$  with  $k_{31}$  knots and the second margin is a thin-plate spline of  $\{\mathbf{e}_i, \mathbf{n}_i\}$  with rank  $k_{32}$ ;
- $f_4$  is a natural cubic spline of  $\mathbf{h}_i$  with  $k_4$  knots;
- $\epsilon_{iy}$  are i.i.d. model errors (the i.i.d. assumption is reasonable as the autocorrelated errors in daily logBS will be aggregated to zero over a year long span).





**Figure 3.26:** [Illustration of annual mean logBS dataset (part 2)] Panel (a) is a boxplot-like plot (see §2.4 if you need an explanation) for logBS from all stations against year. Panel (b) is a boxplot of mean logBS over all years and all stations of a particular environmental type. It is obvious that rural areas had lower logBS. The rest of the boxplot-like graphs sketch logBS from all stations of a particular environmental type against year. logBS declined at different rate at different types of areas, and the difference between rural areas and other areas is the most notable.



**Figure 3.27:** Choosing  $k$  for model 3.4 fails. The  $edf$  for  $\hat{f}_1$  and  $\hat{f}_2$  grows linearly as  $k$  increases.

**Table 3.7:** Refitting model 3.4 by simultaneously increasing  $k$  values for splines and check how their resulting effective degree of freedom change. Again, they grow linearly with  $k$ .

	$\hat{f}_1$	$\hat{f}_2$	$\hat{f}_3$	$\hat{f}_4$	$\tilde{f}_1$
$(k_1, k_2, \{k_{31}, k_{32}\}, k_4)$	$edf$				
$(20, 200, \{20, 200\}, 10)$	15.53	173.16	874.81	7.77	88.31
$(25, 250, \{25, 250\}, 15)$	21.29	211.76	1046.94	11.23	104.74
$(30, 300, \{30, 300\}, 20)$	26.16	251.78	1207.95	16.01	129.22
$(35, 350, \{35, 350\}, 25)$	31.24	287.81	1390.25	21.44	151.93
$(40, 400, \{40, 400\}, 30)$	35.86	328.67	1533.91	25.07	177.60

As usual, practical model development starts from choosing  $k$  values. To start with, consider an initial model fitting with  $k_1 = k_{31} = 20$ ,  $k_2 = k_{32} = 200$  and  $k_4 = 10$ . Extract the partial residuals w.r.t.  $\hat{f}_1(y; 20)$  and fit those partial residuals over  $f_1(y; k_1)$  for increasingly bigger  $k_1$  and see when the effective degree of freedom plateaus. However, at this very first step some problem has occurred: the  $edf$  never plateaus! Maybe we can start with another model component instead? But trying  $f_2(\{e_i, n_i\}; k_2)$  exposes the same issue: the  $edf$  grows linearly with  $k$ . Figure 3.27 demonstrates these results for  $\hat{f}_1$ ,  $\hat{f}_2$  and  $\hat{f}_3$ .

Apart from using partial residuals to iteratively check adequacy of  $k$  values, a more straightforward (but also more computationally expensive) way is to increase all  $k$  values, refit the model and check  $edf$ . Table 3.7 presents the result, and it is concerning that  $edf$  still grow linearly with  $k$ .

### 3.4.3 Including i.i.d. random effects

Observations in the last section imply that there are more variability between data across stations than what smooth functions assume: more variability between mean of logBS across stations than what  $f_2(\{e_i, n_i\}; k_2)$  assumes and more variability between the shape of trend across stations than  $f_1(y; k_1)$  and  $f_3(y, \{e_i, n_i\}; k_{31}, k_{32})$  assume. So when splines are parametrized with increasingly more knots to be more flexible, they continue to approximate data values producing smaller residuals. A question may be raised on whether all between-station difference should be explicitly modelled as mean via complicated regression functions; it may be a good idea to model some between-station difference as random effects and quantify their uncertainty; in this way we focus on modelling a more general and less complicated space-time relationship via spline functions.

Let us thus consider the following model

**Model 3.5:** model 3.4 + i.i.d. random effects

$$\begin{aligned} \log BS_{iy} = & f_0(E_i; 5) + f_1(y; k_1) + \tilde{f}_1(E_i, y; 5, k_1) + \\ & f_2(\{e_i, n_i\}; k_2) + f_3(y, \{e_i, n_i\}; k_{31}, k_{32}) + f_4(h_i; k_4) + \\ & f_5(i; 2626) + f_6(y; 43) + \epsilon_{iy}, \end{aligned}$$

where new components compared with model 3.5 are

**Table 3.8:** Refitting model 3.5 by simultaneously increasing  $k$  values for splines and check how their resulting effective degree of freedom change.

	$\hat{f}_6$	$\hat{f}_5$	$\hat{f}_1$	$\hat{f}_2$	$\hat{f}_3$	$\hat{f}_4$	$\tilde{f}_1$
$(k_1, k_2, \{k_{31}, k_{32}\}, k_4)$	<i>edf</i>						
(20, 200, {20, 200}, 10)	35.07	2231.29	5.00	120.88	1298.83	1.00	25.94
(25, 250, {25, 250}, 15)	35.04	2194.11	5.02	140.90	1557.71	1.00	25.95
(30, 300, {30, 300}, 20)	35.07	2162.51	5.03	153.03	1868.25	1.00	25.97
(35, 350, {35, 350}, 25)	35.05	2137.41	5.04	162.86	2114.05	1.00	25.99
(40, 400, {40, 400}, 30)	35.10	2117.62	5.02	174.54	2363.10	1.00	26.00

- $f_5$ , an i.i.d. random effect for station  $i$ ;
- $f_6$ , an i.i.d. random effect for year  $y$ .

The random effect for year should be interpreted as an adjustment to the station random effect, so that the random effect for station  $i$  in year  $y$  is  $f_5(i) + f_6(y)$ . Estimation of random effects can be easily incorporated into standard estimation procedure of addition models with penalized regression. The design matrix of a random effect is just an matrix with dummy columns, and the associated penalty matrix is an identity matrix giving a ridge penalty. Given an estimated smoothing parameter  $\hat{\lambda}_5$  for  $\hat{f}_5$  for example, as well as an estimated residual variance  $\hat{\sigma}_\epsilon^2$ , the variance of  $\hat{f}_5$  is just  $\hat{\sigma}_\epsilon^2 / \hat{\lambda}_5$ . Predicted  $\hat{f}_5$  is obtained straightforward in the solution of penalized least squares.

Because of the existence of random effects,  $k$  values can not be reliably chosen using partial residuals. There is a trade-off between the complexity in  $f_2$  and  $f_5$  (and likewise between  $f_1$  and  $f_6$ ) that can only be estimated jointly during model fitting. So let us take the most expensive approach: increase  $k$  values and refit the whole model. The result is presented in Table 3.8, and the following observations can be made.

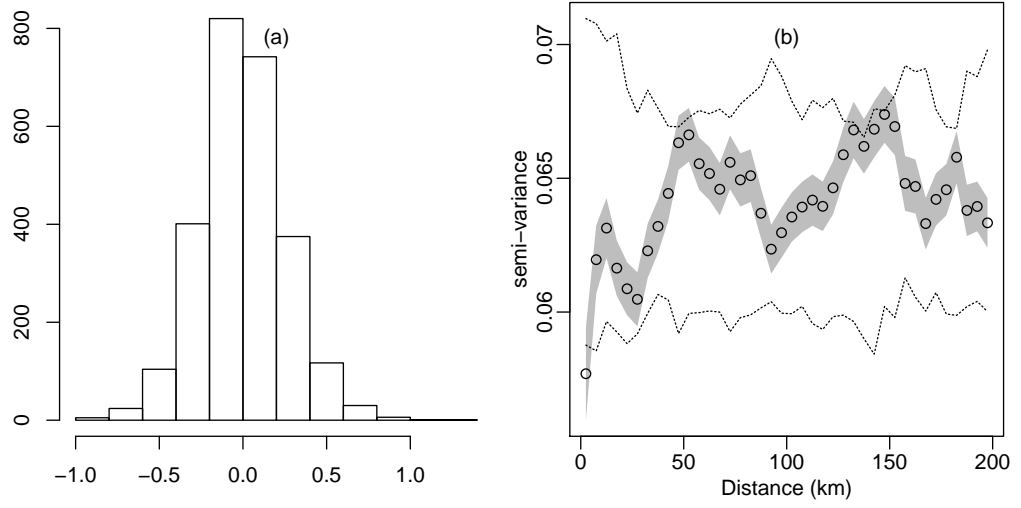
- With the introduction of random effects, *edf* of  $\hat{f}_1$ ,  $\hat{f}_4$  and  $\tilde{f}_1$  plateaus respectively at about 5, 1 and 26.
- The *edf* of  $\hat{f}_2$  does not strictly plateau, but it grows sublinearly with  $k_2$ . For example, for *edf* :  $k_2$  ratio is 0.6 at  $k_2 = 200$ , 0.56 at  $k_2 = 250$ , 0.51 at  $k_2 = 300$ , etc., which keeps dropping as  $k_2$  increases linearly. (As a comparison, this ratio is steadily about 0.8 in Table 3.7.) In fact, slightly less than half of the coefficients are already suppressed at  $k_2 = 200$ , so this  $k$  value is readily sufficient.

### 3.4.4 Model summary and checking

$k_1 = k_{31} = 20$ ,  $k_2 = k_{32} = 200$  and  $k_4 = 10$  are sufficient for model 3.5, so let us now summarize the model and do some checking.

The effective degree of freedom for all spline functions have been given in Table 3.8. In total 6750 parameters have been used for  $n = 24239$  data and the fitted model has an *edf* of 3730.7. The estimated residual variance is 0.02939 (residual sum of squares divided by  $(n - \text{edf})$ ); the variance of logBS data is 0.75622, so the fitted model has an adjusted R-squared at  $1 - 0.02939 / 0.75622 = 96.11\%$  which is very high. The estimated variance for station specific random effects  $\hat{f}_5$  and year specific random effects  $\hat{f}_6$  are respectively 0.0741 (much larger than residual variance) and 0.00272 (much smaller than residual variance).

One concern is whether the predicted random effects are really i.i.d.. If the station specific random effects turn out to exhibit spatial autocorrelation, or the year specific random effects show temporal



**Figure 3.28:** Predicted station-specific random effects  $\hat{f}_5$  in fitted model 3.5. Panel (a) is their histogram, and panel (b) is their empirical variogram (computed using `geoR` package). The dots in the graph are the classic (i.e., moment) estimates for variogram, and the gray ribbon covers region within 2 standard deviation of the estimate. Envelops (dashed lines) for the estimates (by permutation of the data values on the spatial locations) are added for testing the existence of spatial autocorrelation. No violation of the i.i.d. assumption for the random effects is seen.

autocorrelation, then it is more appropriate to model them with some covariance or correlation matrix. Panel (b) of Figure 3.28 is an empirical variogram of predicted  $\hat{f}_5$ , and no spatial correlation is spotted. The ACF in panel (b) of Figure 3.29 confirms that there is no temporal autocorrelation between predicted  $\hat{f}_6$ .

Standard residual checking involves inspecting residuals against fitted values and independent variables, as is done in Figure 3.30. No violation of constant error variance assumption is seen from panel (a), and no unmodelled trend is spotted in the rest of panels. The latter is not surprising, since all these independent variables have been included in the model and sufficient  $k$  values have been chosen for all splines.

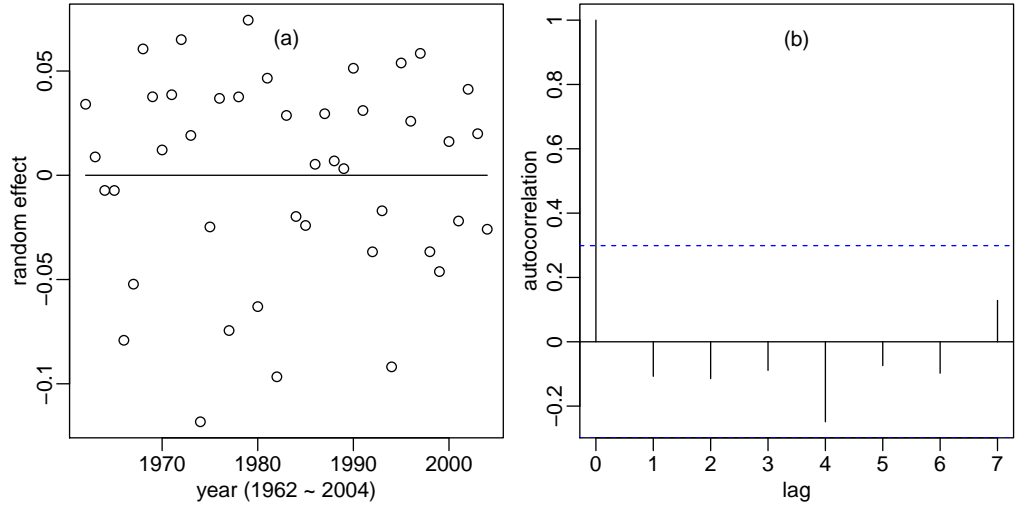
In spatial and spatial-temporal modelling it is also important to check if there are unmodelled spatial autocorrelation in residuals. Figure 3.31 produces empirical variograms for spatial residuals every three years. It seems that the thin-plate spline  $f_2(\{\mathbf{e}_i, \mathbf{n}_i\}; k_2)$  and the tensor product spline  $f_3(\mathbf{y}, \{\mathbf{e}_i, \mathbf{n}_i\}; k_{31}, k_{32})$  have adequately modelled all spatial-temporal variability.

### 3.4.5 Visualization of model prediction

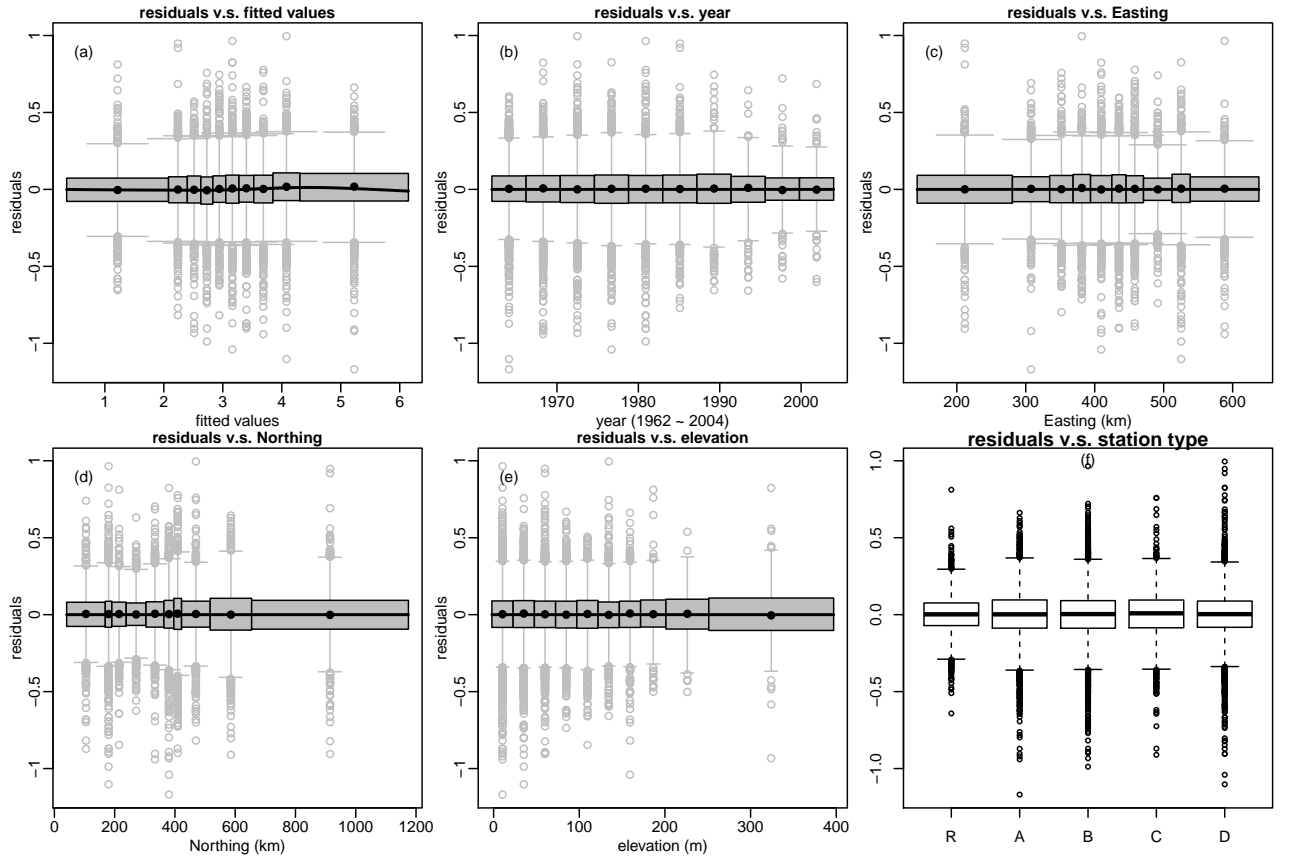
So far the fitted model 3.5 looks reasonable. Now let us visualize its components and demonstrate how it makes prediction.

Figure 3.32 illustrates various main effects in the fitted 3.5.

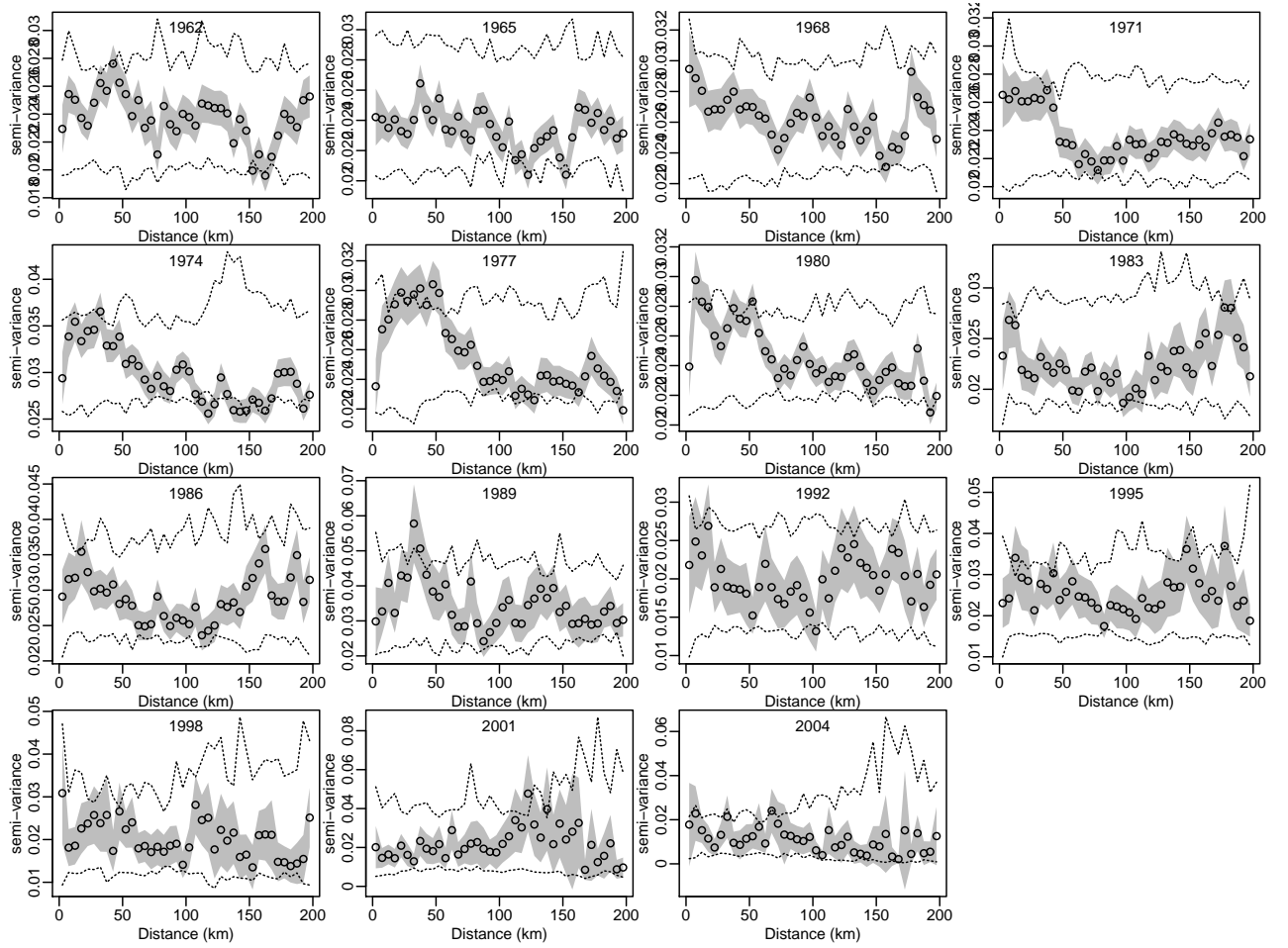
- Panel (a) shows the estimated mean long term trend  $\hat{f}_1$ . Its complexity appears to be somewhere between a linear line and a quadratic polynomial. Or it looks like two linear line segments that join at around year 1980. Before 1980 the lines descends fast, and slower afterwards.
- Panel (b) is the estimated linear line for elevation effect  $\hat{f}_4$ . That the confidence interval shrinks to zero somewhere at the middle may be confusing, but it is expected since in additive models all additive components are subject to centring constraints for better identifiability (and that is why all splines in the Figure cross zero). A linear line has only 1 degree of freedom, thus has no degree of freedom at centring points.



**Figure 3.29:** Predicted year-specific random effects  $\hat{f}_6$  in fitted model 3.5. Panel (a) is a scatter plot of  $\hat{f}_6$  against year. A smoothing spline (black line) is fitted to the data but ends up with a horizontal line at 0, confirming that there is no trend in  $\hat{f}_6$ . Panel (b) is an autocorrelation function (ACF) for  $\hat{f}_6$ . Although these random effects are ordered in time, no autocorrelation is spotted so the i.i.d. assumption is validated.



**Figure 3.30:** Inspecting residuals of fitted model 3.5 against fitted values (panel (a)) and other covariates ((b) to (f)). No violation of constant error variance assumption is seen from panel (a), and no unmodelled trend is spotted in the rest of panels. The latter is not surprising, since all these independent variables have been included in the model and sufficient  $k$  values have been chosen for spline functions. See §2.4 if you need an explanation for these boxplot-like graphs.

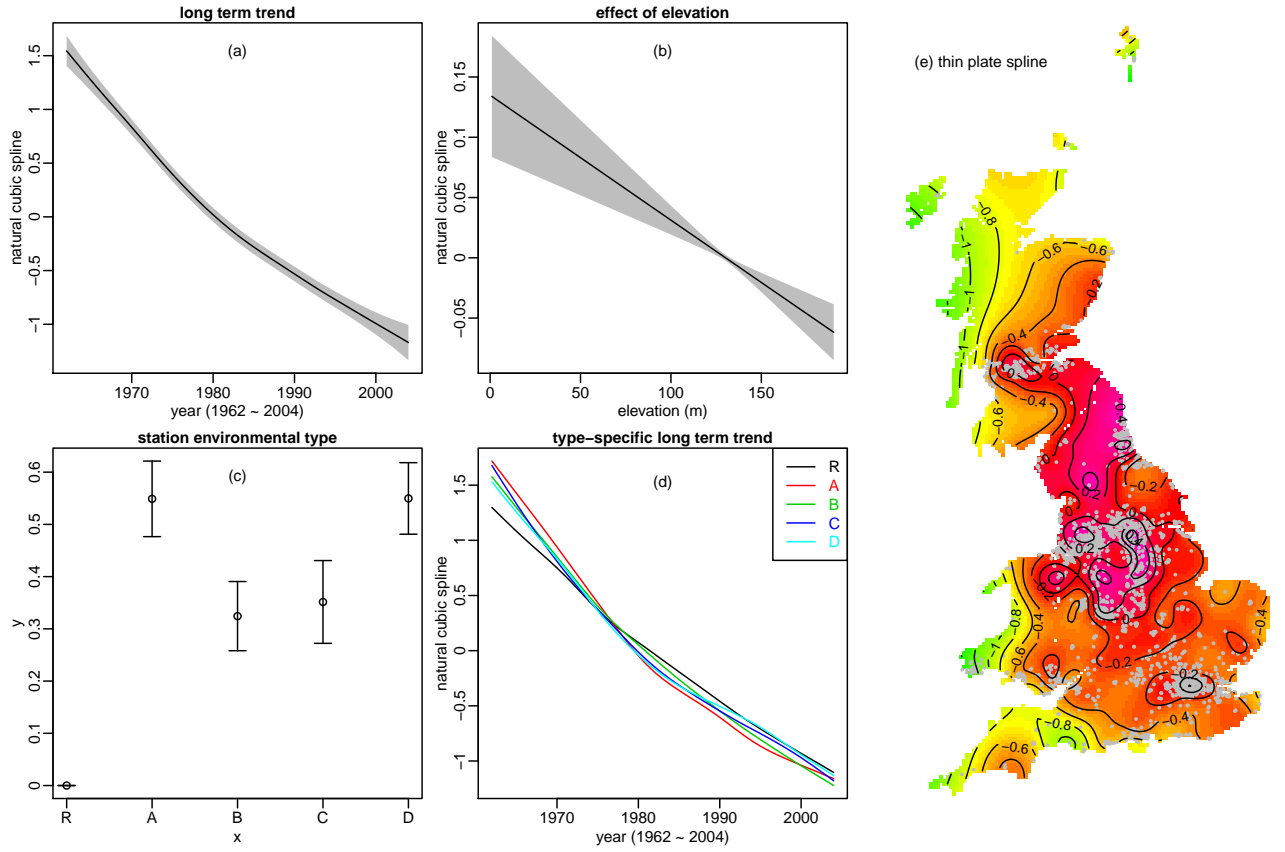


**Figure 3.31:** Residual checking of fitted model 3.5 in spatial domain every three years, i.e., at 1962, 1965, ..., 2004. No unmodelled spatial autocorrelation is seen.

- Panel (c) is the estimated mean logBS at different environmental areas. “Rural” (“R”) is the reference level for this factor variable hence is contrasted, i.e., it has zero mean and zero standard error. Estimates at other levels are in fact the mean difference from “R”. It is clear that all other areas have higher mean than rural areas. This graph verifies what was previously observed from data (see the boxplot in panel (b) of Figure 3.26).
- Panel (d) is the estimated mean long term trend  $\hat{f}_1 + \tilde{f}_1$  for annual logBS at different types of areas / stations. I choose to illustrate this quantity because it is more interpretable than  $\tilde{f}_1$  alone that describes the mean trend difference between a level and the reference level “R”. It is clear that logBS declined at a different rate in rural areas; the decline rate for other four levels are, while different, more similar. This estimated result agrees with what was observed from data (see Figure 3.26).
- Panel (e) is the estimated thin-plate spline  $\hat{f}_2$  for time invariant mean logBS over space. A contour plot is used for surface visualization as it is more informative / readable than a perspective plot. Pixels are coloured for further enhancement. The colouring scheme used is (roughly) green, yellow, orange, red, purple from low logBS to high logBS values. It can be seen that central England and Northeastern England are areas with high Smoke concentration.

Figure 3.33 illustrates how the fitted model 3.5 predicts logBS when its model terms are added together. 4 stations with different number of operation years are taken as examples. See Figure caption for more details.

The visualization of tensor product spline  $\hat{f}_3$  is much more difficult as it is a 4D surface. Treating the 2D spatial coordinates as a single dimension allows us to inspect the transects line along the time dimension, that is, how the long term trend at a specific spatial location differs from the mean trend

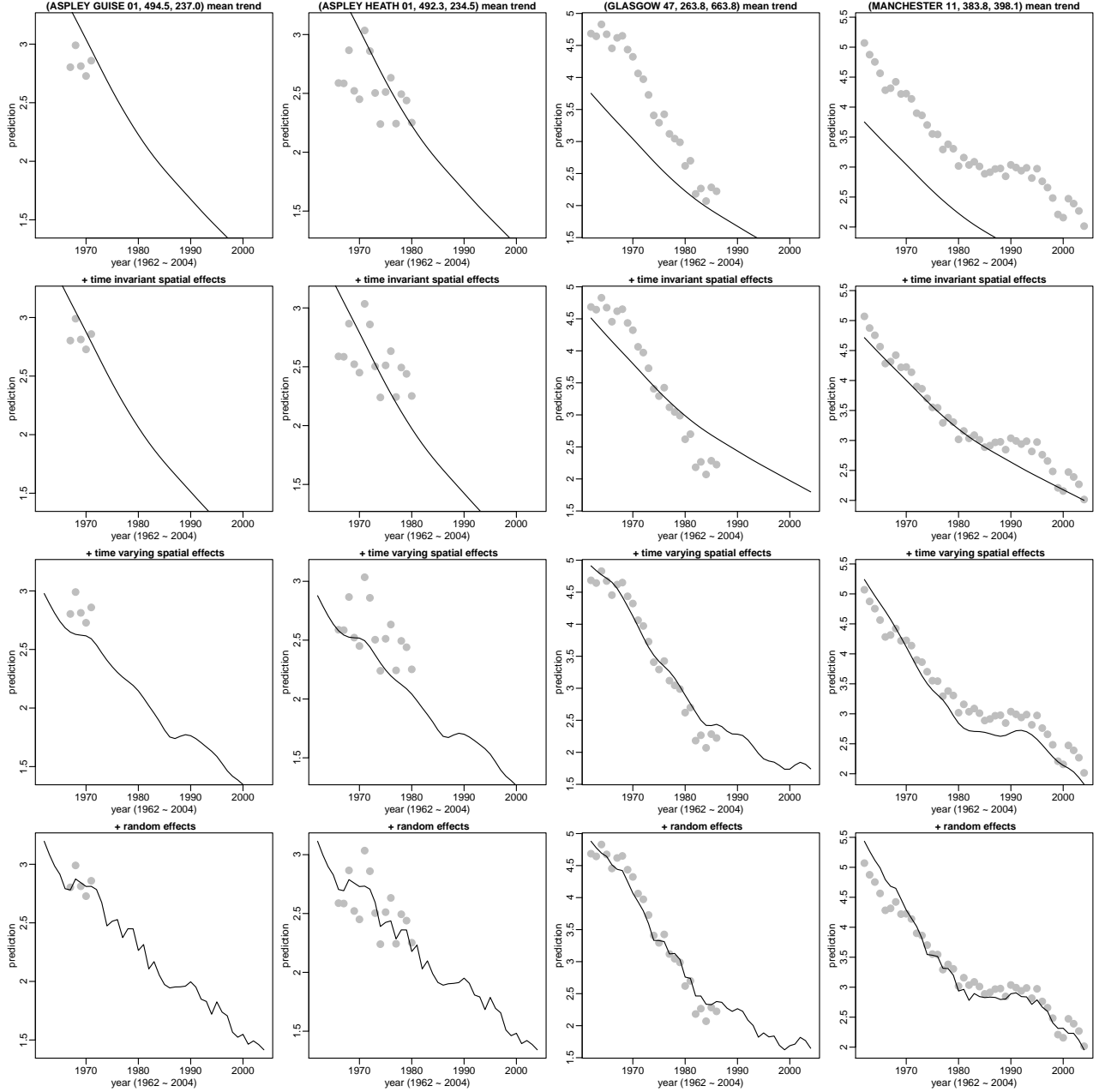


**Figure 3.32:** Estimated main effects in fitted model 3.5. Panel (a) is the estimated mean long term trend  $\hat{f}_1$  for annual logBS (the gray ribbon covers regions within 2 standard error from the estimate); panel (b) is the estimated spline / line  $\hat{f}_4$  for elevation; panel (c) is the estimated mean logBS  $\hat{f}_0$  at different types of areas / stations (the whiskers stretch to 2 standard error from the estimate; note that “rural” type is the reference level and is contrasted, so it has mean 0 and standard error 0); panel (d) is the estimated mean long term trend  $\hat{f}_1 + \tilde{f}_1$  for annual logBS at different types of areas / stations; panel (e) is the estimated thin-plate spline  $\hat{f}_2$  for time invariant mean logBS over space. See main text for more detailed explanations.

$\hat{f}_1$ . Given the simple smooth function in Figure 3.32, the high variability in data (see Figure 3.26) but high R-squared of the fitted model, one can expect this tensor product spline to be complicated. Figure 3.34 sketches transects of the spline through all 2626 stations. Generally those transects are multimodal “wave” like curves so that the mean trend  $\hat{f}_1$  can be “twisted” when added by them. However, the four evident fast ascending lines are very concerning, as they imply that logBS will be predicted to soar at associated stations. This is very likely due to the extrapolation of splines when for example, a station was closed at early years but prediction is asked for the entire 43 years. With a bit of efforts, it can be identified that those four lines are for the four stations in Shetland which are far apart from the majority stations in the Network. (Referring to the map in Figure 3.25, they are on the northernmost small island.) Figure 3.35 demonstrates the extrapolation at those stations.

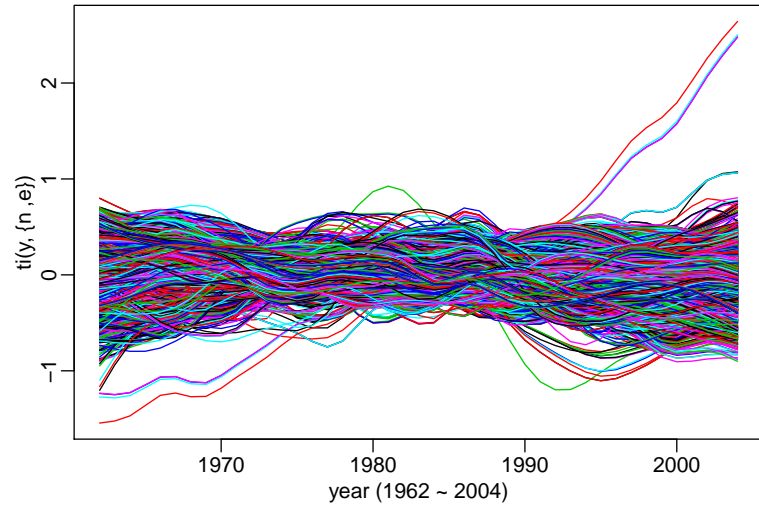
It is then natural to consider removing those stations from Shetland when fitting model, however, the extrapolation problem exists with other stations, too. Some more sophisticated solution is needed to deal with this problem. For the moment, I will set aside this issue.

The aim of spatial-temporal modelling is to make prediction away from sampling locations, or simply put, to produce high-resolution mapping of logBS over space and time. Figure 3.36 illustrates this, by predicting logBS on a 5km  $\times$  5km grid over Great Britain every 8 years. The general decline of logBS can be clearly seen. Note that the extrapolation problem for Shetland can also be spotted.

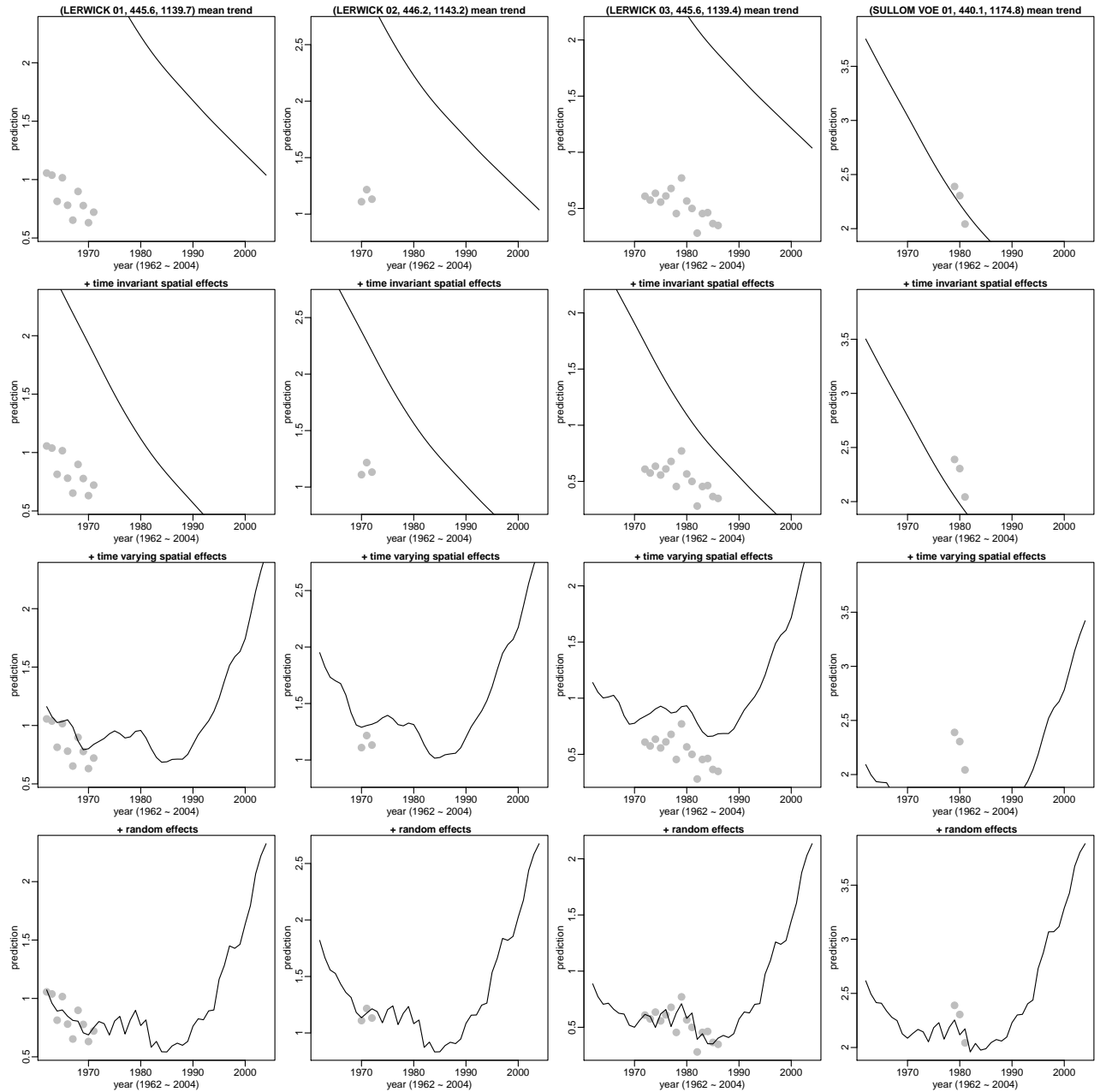


**Figure 3.33:** An illustration of how the fitted model 3.5 predicts logBS when its model terms are added together. Four stations with different length of monitoring history are chosen. The first (leftmost) column panel is for station *ASPLEY GUISE 01* at {494.5 (Easting), 237.0 (Northing)}; it operated for 5 years. The second column panel is for station *ASPLEY HEATH 01* at {492.3, 234.5}, which operated for 15 years. Note that this station is located very close to the first one. The third column panel is for station *GLASGOW 47* at {263.8, 663.8}, with 25 years' service. The final column panel is for the *Manchester 11* that has been thoroughly studied in a previous section. For column panels, from top rows to bottom rows, these model components are added: 1) Mean long term trend  $\hat{f}_1$  (model intercept included). This component is identical for all stations, although the curves in the graphs look different due to different scale of the vertical axis; 2) Time invariant spatial effects  $\hat{f}_2 + \hat{f}_4 + \hat{f}_0$ ; This component yields a vertical shift to the previous component; 3) time varying spatial effects  $\hat{f}_3 + \hat{f}_1$ ; This component “twists” the curve so that it has a closer shape to that of data. 4) predicted random effects  $\hat{f}_5 + \hat{f}_6$ . This component adds another vertical shift as well as fine-scaled adjustment (that is not easily modelled by a smooth function) to the predicted curve. As more terms are added, the predicted curve approximates data better. Note that the final curves for the first two stations look extremely similar. This is because the stations are very closely located so there is little spatial variability in prediction.

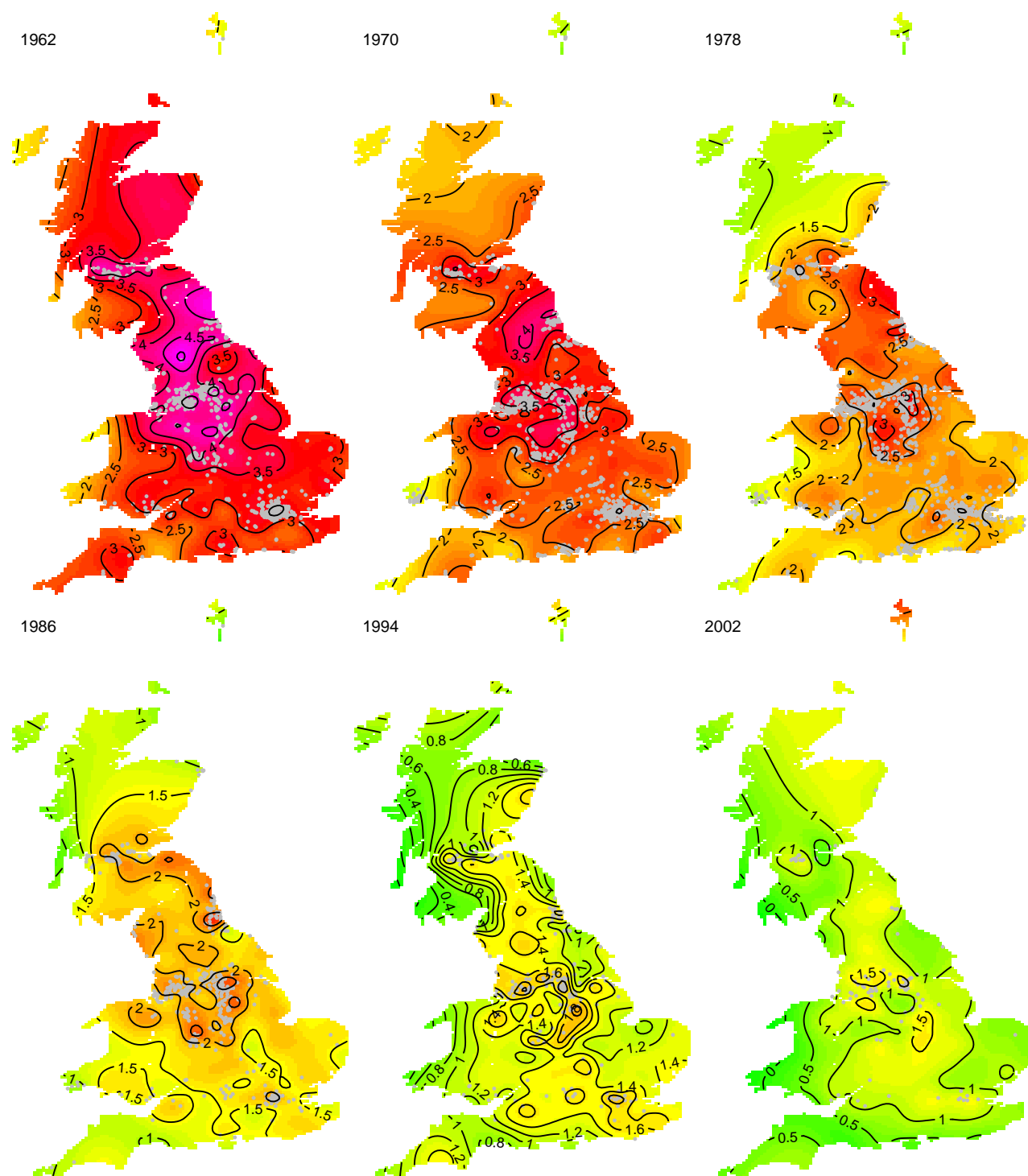




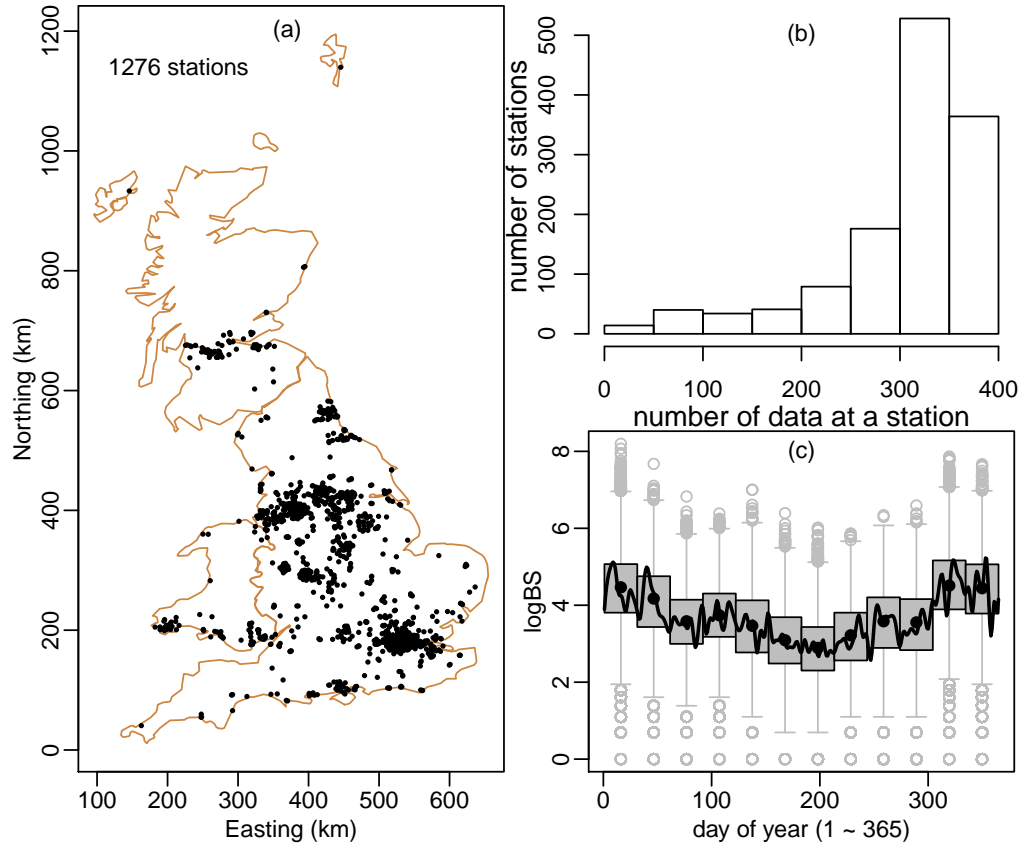
**Figure 3.34:** Transects of tensor product spline  $\hat{f}_3$  at all stations. Most of the curves are multimodal “wave” like curves which only lead to mild “twist” of  $\hat{f}_1$ , but there are four evident fast ascending lines associated with four stations at Shetland that would cause unreasonable extrapolation when predicting logBS. See Figure 3.35.



**Figure 3.35:** The extrapolation problem at Shetland, Scotland.



**Figure 3.36:** Spatial-temporal prediction of logBS on a 5km  $\times$  5km spatial grid over Great Britain every 8 years using fitted model 3.5.



**Figure 3.37:** [Illustration of daily logBS dataset in year 1967 (part 1)] Panel (a) shows all 1276 operating stations in year 1967. Panel (b) is a histogram for the number of stations with a certain number of daily data. The majority of stations had operated for over 300 days. Panel (c) is a boxplot-alike graph for logBS from all stations against day of year.

### 3.4.6 Summary

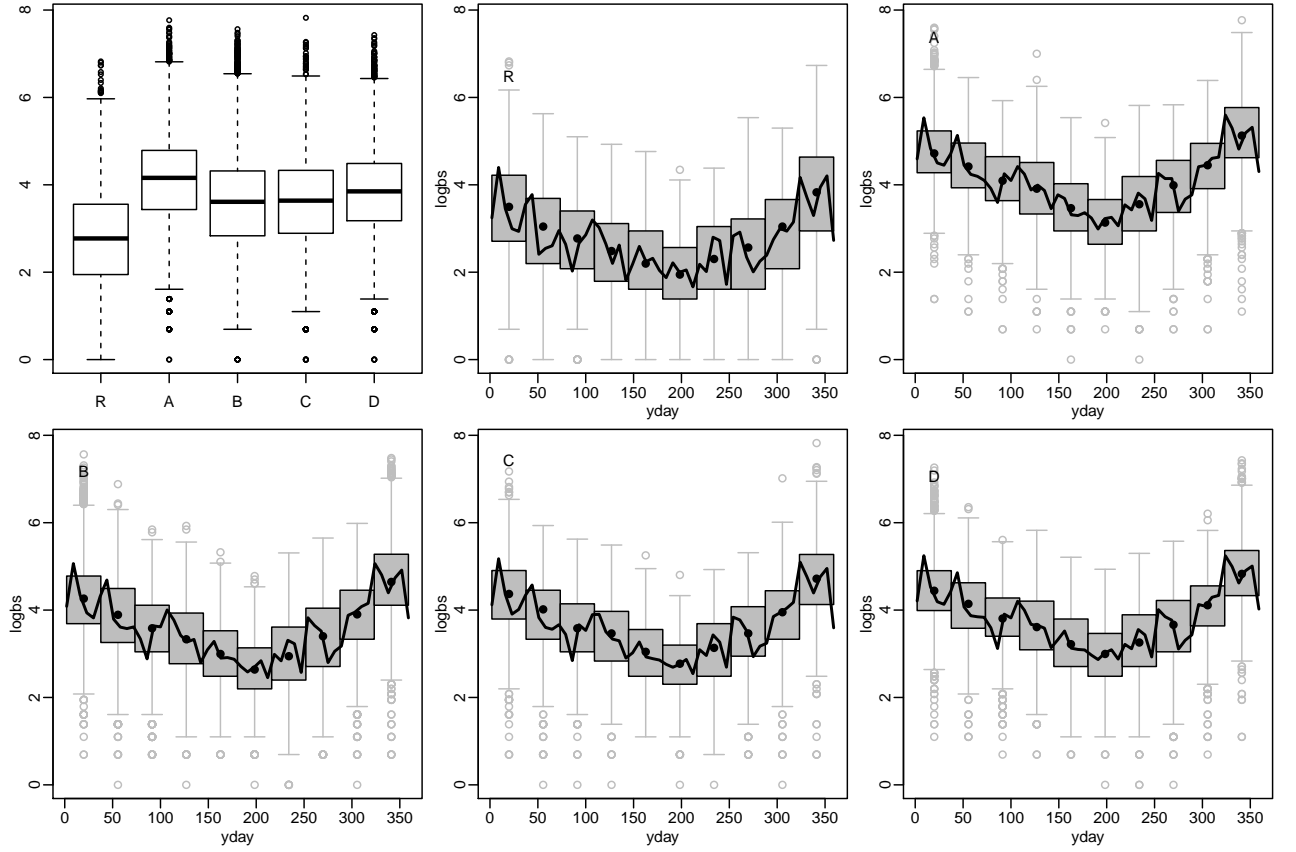
In this section the annual mean logBS dataset is used to demonstrate how additive models with penalized regression splines can be used for spatial-temporal modelling. The following message can be learned.

- Spline functions alone can not adequately model all spatial or temporal variability in logBS data. I.i.d. random effects on both space and time domain are needed to model non-smooth change if logBS over space and time.
- The extrapolation problem is inherent for splines. Some further investigation on this issue is required in future.

## 3.5 Spatial-temporal modelling for daily logBS in 1967

### 3.5.1 Description of data

In this section, I will attempt developing a model for daily logBS in year 1967, when the size of the Network was at its peak with 1276 sites. The dataset comprises 387253 daily data, approximately 4.1% of the complete daily logBS dataset. See Figure 3.37 and Figure 3.38 for some visualization of this dataset.



**Figure 3.38:** [Illustration of daily logBS dataset in year 1967 (part 3)] The first boxplot is for the mean logBS over all days and all stations of a particular environmental type. It is clear that logBS differs in mean at different environmental areas. The subsequent boxplot-like graphs sketch logBS against day of year for all stations of the same environmental type. There does not seem to be any difference between the mean seasonality of logBS at different environmental areas.

### 3.5.2 Building a model for Monday data

Since daily data has strong autocorrelation, I will thin the data and first investigate Monday data. This subset of data has 54386 daily data from 1275 stations (1 fewer than the total number stations in the year).

In fact, the modelling experience established in previous sections has provided rich information to help building a model here.

- From Figure 3.24 in §3.3, it can be learned that seasonality of logBS and temperature effects on logBS vary spatially, so tensor product splines are needed to model interaction between spatial coordinates  $\{\mathbf{e}_i, \mathbf{n}_i\}$  and week of year  $\mathbf{w}_y$ , as well as  $\{\mathbf{e}_i, \mathbf{n}_i\}$  and temperatures  $T_{id}^0, T_{id}^*$ .
- From §3.4.6, it can be learned that station-specific random effects and week-specific random effects are necessary to model non-smooth change in spatial and temporal variability that is difficult to model by spline functions.

In addition,

- Figure 3.38 shows no noticeable difference in (the shape of) seasonality of logBS between different environmental areas, so there is no need to model interaction between week of year  $\mathbf{w}_y$  and environmental types  $\mathbf{E}_i$ .

- The number of knots of natural cubic splines or the rank of thin-plates that are adequate for modelling the effects of various variables are also known.

So it is straightforward to propose the following model as a start:

**Model 3.6:** A model for Monday logBS

$$\begin{aligned} \log\text{BS}_{i\mathbf{w}} = & f_0(\mathbf{E}_i) + f_1(\mathbf{w}_y; 20) + f_2(\{\mathbf{e}_i, \mathbf{n}_i\}; 200) + f_3(\mathbf{w}_y, \{\mathbf{e}_i, \mathbf{n}_i\}; 20, 200) + \\ & f_4(\mathbf{T}_{i\mathbf{d}}^0; 15) + f_5(\mathbf{T}_{i\mathbf{d}}^0, \{\mathbf{e}_i, \mathbf{n}_i\}; 15, 200) + \\ & f_6(\mathbf{T}_{i\mathbf{d}}^*; 15) + f_7(\mathbf{T}_{i\mathbf{d}}^*, \{\mathbf{e}_i, \mathbf{n}_i\}; 15, 200) + \\ & f_8(\mathbf{T}_{i\mathbf{d}}^0, \mathbf{T}_{i\mathbf{d}}^*; 15, 15) + f_9(i; 1275) + f_{10}(\mathbf{w}_y; 52) + f_{11}(\mathbf{h}_i; 10) + \epsilon_{i\mathbf{w}}, \end{aligned}$$

where

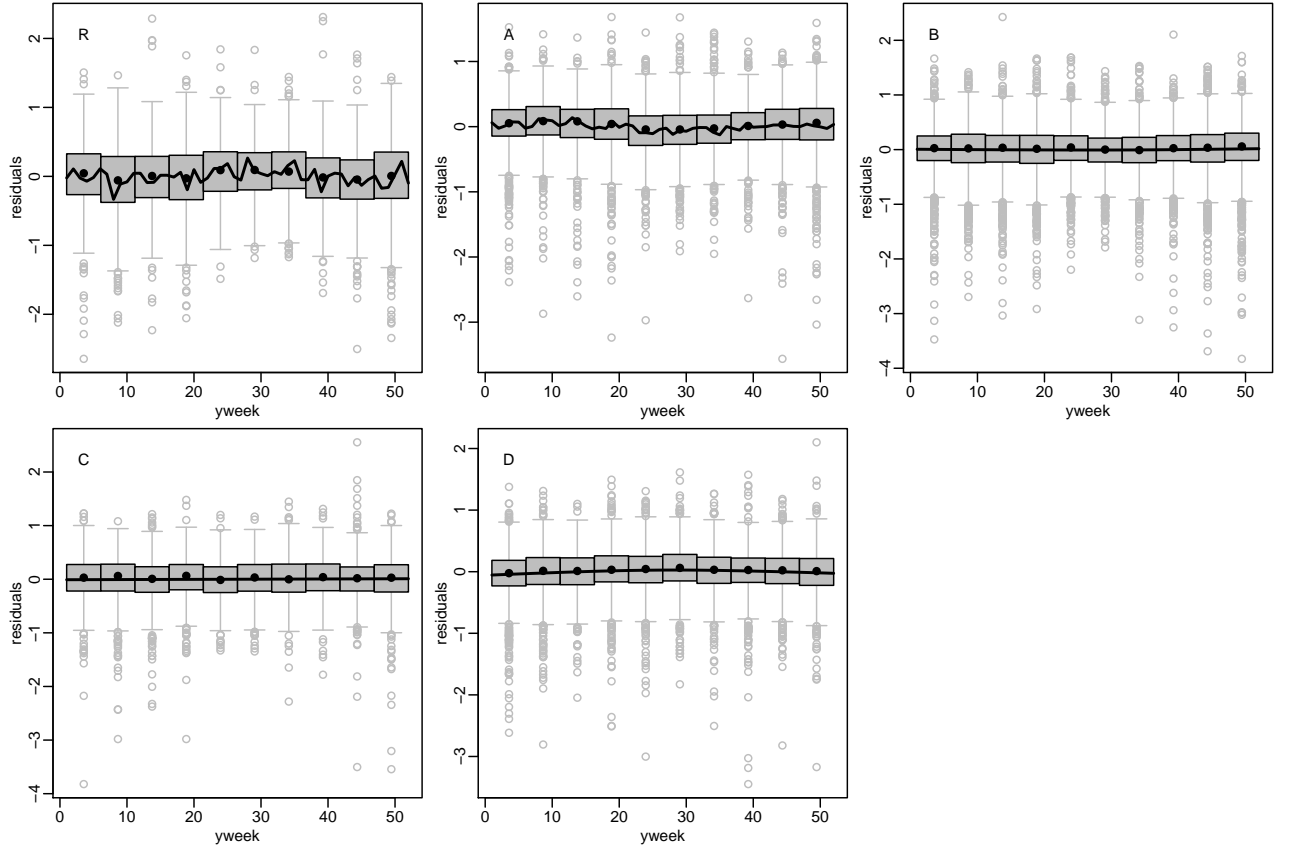
- $f_0(\mathbf{E}_i; 5)$  is a random effect modelling the mean of logBS from stations of environmental type  $\mathbf{E}_i$ ;
- $f_1(\mathbf{w}_y; 20)$  is a cubic cyclic spline with 20 knots for  $\mathbf{w}_y$ , modelling the mean seasonality of logBS at all stations;
- $f_2(\{\mathbf{e}_i, \mathbf{n}_i\}; 200)$  is a rank-200 thin-plate spline for spatial coordinates of stations, modelling the mean logBS over space;
- $f_3$  is a tensor product spline modelling how station-specific seasonality deviates from  $f_1$ ; its first margin is a cubic cyclic spline with 20 knots for  $\mathbf{w}_y$  and its second margin is a rank-200 thin-plate spline for spatial coordinates of stations;
- $f_4(\mathbf{T}_{i\mathbf{d}}^0; 15)$  is a natural cubic spline with 15 knots for daily minimum temperature  $\mathbf{T}_{i\mathbf{d}}^0$ , modelling the mean effect of this variable at all stations;
- $f_5$  is a tensor product spline modelling how station-specific  $\mathbf{T}_{i\mathbf{d}}^0$  effect deviates from  $f_4$ ; its first margin is a natural cubic spline with 15 knots for  $\mathbf{T}_{i\mathbf{d}}^0$  and its second margin is a rank-200 thin-plate spline for spatial coordinates of stations;
- $f_6(\mathbf{T}_{i\mathbf{d}}^*; 15)$  is a natural cubic spline with 15 knots for diurnal temperature difference  $\mathbf{T}_{i\mathbf{d}}^*$ , modelling the mean effect of this variable at all stations;
- $f_7$  is a tensor product spline modelling how station-specific  $\mathbf{T}_{i\mathbf{d}}^*$  effect deviates from  $f_6$ ; its first margin is a natural cubic spline with 15 knots for  $\mathbf{T}_{i\mathbf{d}}^*$  and its second margin is a rank-200 thin-plate spline for spatial coordinates of stations;
- $f_8(\mathbf{T}_{i\mathbf{d}}^0, \mathbf{T}_{i\mathbf{d}}^*; 15, 15)$  is a tensor product spline modelling the interaction between two temperature variables;
- $f_9(i; 1275)$  is an i.i.d. random effect for 1275 stations;
- $f_{10}(\mathbf{w}_y; 52)$  is an i.i.d. random effect for week of year (there are 52 Mondays in this year);
- $f_{11}(\mathbf{h}_i; 10)$  is a cubic spline for elevation  $\mathbf{h}_i$ ;
- $\epsilon_{i\mathbf{w}}$  is an i.i.d. model error.

### 3.5.3 Model summary and checking

Table 3.9 summarizes the effective degree of freedom for all splines in fitted model 3.6. They are much smaller than the specified  $k$  values so  $k$  values are adequate. This model involves 10932 parameters with an *edf* of 3928.32 for 54386 data. The estimated residual variance is 0.1954, and the sample

**Table 3.9:** effective degree of freedom for all splines in fitted model 3.6

$\hat{f}_0$	$\hat{f}_1$	$\hat{f}_2$	$\hat{f}_3$	$\hat{f}_4$	$\hat{f}_5$	$\hat{f}_6$	$\hat{f}_7$	$\hat{f}_8$	$\hat{f}_9$	$\hat{f}_{10}$	$\hat{f}_{11}$
3.98	4.74	106.16	1419.05	4.81	624.35	4.50	520.67	88.27	1103.86	45.97	1.00

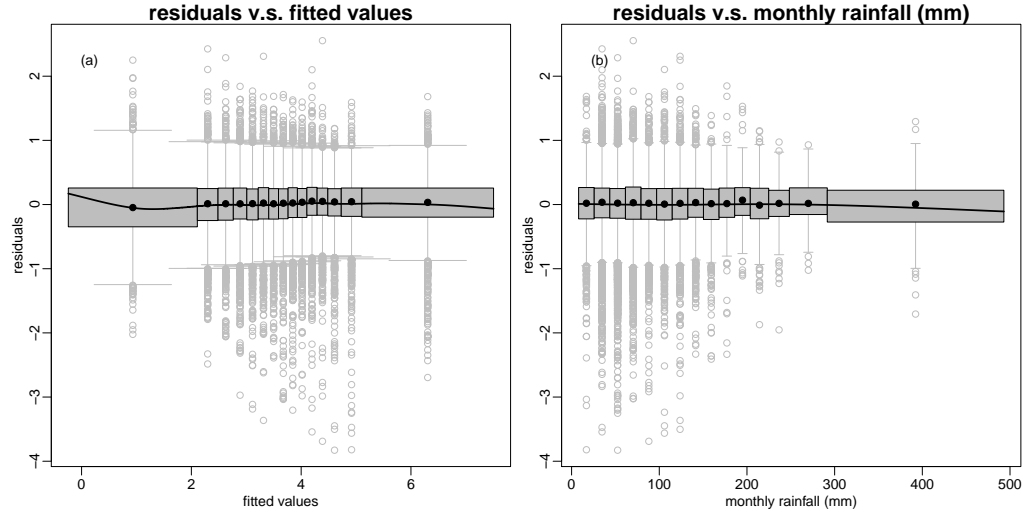


**Figure 3.39:** Inspecting residuals of fitted model 3.6 from each environment area against week of year. For type “B”, “C” and “D”, there is strictly no unmodelled trend; for type “R” and “A”, the smoothing spline has fitted some trend, but it looks more like random noise. So there is no need to include an interaction between type and week of year in the model.

variance for logBS data is 1.2642, so the fitted model has an adjusted R-squared of 84.54%. Estimated variance for weekly random effects and station-specific random effects are respectively 0.0892 and 0.1116.

As I commented earlier on Figure 3.38, there does not seem to be any difference between seasonality of logBS from different environment areas. However, it will still be good to confirm this from residuals of fitted model 3.6. Figure 3.39 produces boxplot-alike graphs for residuals from different environment areas against week of year. For type “B”, “C” and “D”, there is strictly no unmodelled trend; for type “R” and “A”, the smoothing spline has fitted some trend, but it looks more like random noise. So I believe there is no need to include an interaction between type and week of year in the model.

Figure 3.40 does some other checking on residuals. The “residuals v.s. fitted” plot assures that residuals have constant variance independent of regression mean. There is no need to plot residuals against variables included in the model; as is previously demonstrated in Figure 3.30, there will be no unmodelled trend w.r.t. those variables when  $k$  values are sufficient. It is more interesting to plot residuals against rainfall to see if rainfall can potentially improve the model. However, from panel (b) it appears that the effect of rainfall is negligible. Only rainfall over 300 mm per month seems to reduce logBS slightly, however, there are not many rainfall records above this threshold, so such conclusion is not strong. Once again, the effects of rainfall have to be investigated when there are more data. Panel (c) is a boxplot for autocorrelation function (ACF) of residuals from every stations. On average, there is no sign of residual autocorrelation. This is somehow expected, since I have thinned daily data only retaining those from Monday for model development.



**Figure 3.40:** Some residual checking for fitted model 3.6. Panel (a) inspects residuals against fitted values and residuals appear to have a constant variance that is independent of regression mean. Panel (b) plots residuals against rainfall, a variable that is not yet in the model. It appears that the effect of rainfall is negligible. Only rainfall over 300 mm per month seems to reduce logBS slightly, however, there are not many rainfall records above this threshold, so such conclusion is not strong. Panel (c) is a boxplot for autocorrelation function (ACF) of residuals from every stations. On average, there is no sign of residual autocorrelation so the i.i.d. model error assumption is validated.

Random effects must also be checked to validate their i.i.d. assumption. Figure 3.41 includes a histogram and ACF for predicted weekly random effects. No correlation is seen at all lags great than 0. Figure 3.42 includes a histogram and empirical variogram for predicted station-specific random effects. There is no sign of spatial autocorrelation.

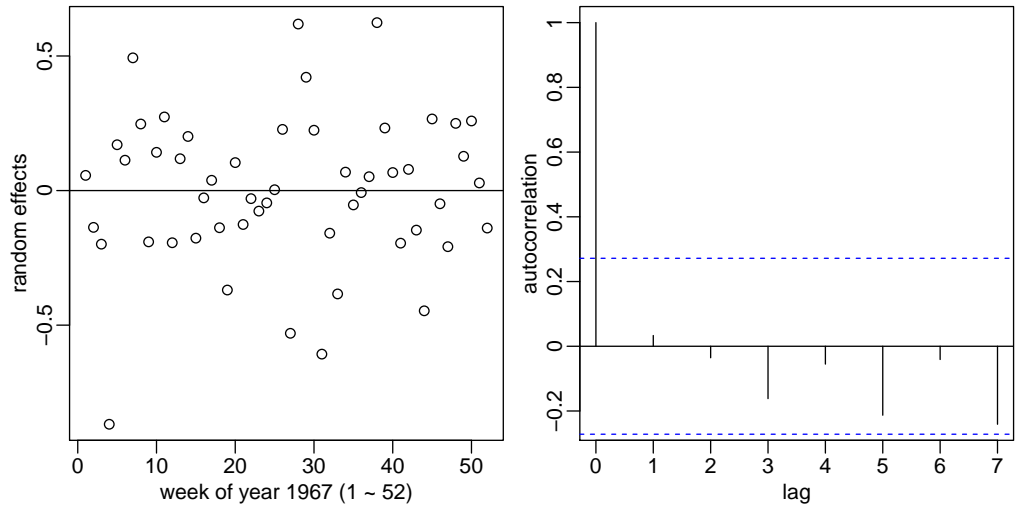
Now let us visualize the estimated splines. Figure 3.43 illustrates main effects in fitted model 3.6, namely  $\hat{f}_1$  (panel (a)),  $\hat{f}_2$  (panel (e)),  $\hat{f}_{11}$  (panel (b)),  $\hat{f}_4$  (panel (c)) and  $\hat{f}_6$  (panel (d)). They all turn out simple; of course, the tensor product splines modelling their spatial variation are more complicated. Transects plot like Figure 3.34 can be produced, but they are of little interest. It is more interesting to visualize  $\hat{f}_8$ , the interaction between two temperature variables. Figure 3.44 visualizes this 3D surface via a perspective plot. The surface is much more complicated compared with the Figure 3.11 from *Manchester 11* case study. This demonstrates the impact of sample size on estimation.

In spatial-temporal modelling residuals should also be check to see if there is unmodelled spatial autocorrelation. Similar to Figure 3.31, I proceed to produce the empirical variograms for spatial residuals in each Monday. Figure 3.45 presents those 52 variograms. On most of the days the fitted model is adequate, leaving no spatial autocorrelation in residuals. But it does leave an excess amount of autocorrelation on day 121.

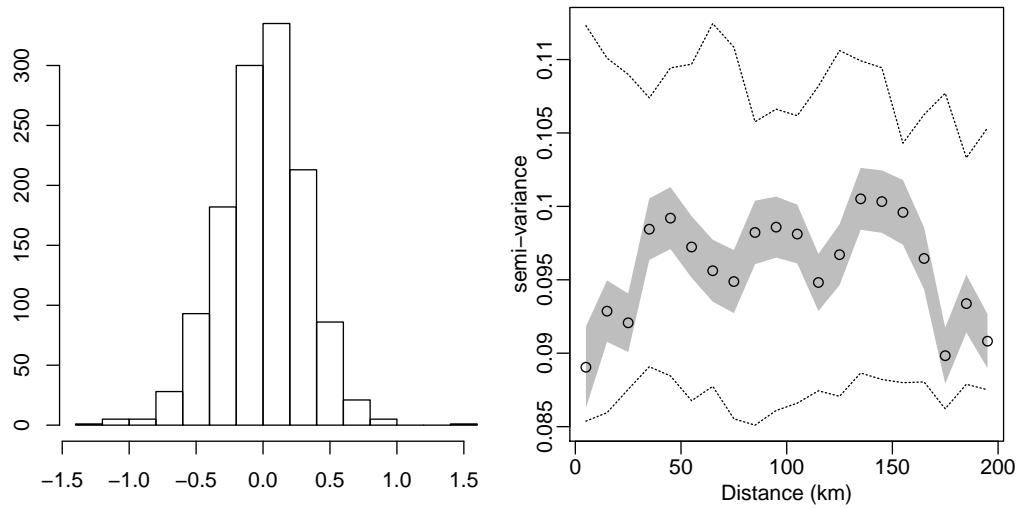
### 3.5.4 Removing spatial autocorrelation in residuals with three-way interaction

It is not impossible to add more components to model 3.6 to remove spatial autocorrelation in its residuals. Since the residuals have unmodelled spatial structure, the new components are likely to involve  $\{\mathbf{e}_i, \mathbf{n}_i\}$ . However, given that the marginal effect  $f_2(\{\mathbf{e}_i, \mathbf{n}_i\})$  and all two-way interactions  $f_3(\mathbf{w}_y, \{\mathbf{e}_i, \mathbf{n}_i\})$ ,  $f_5(\mathbf{T}_{id}^0, \{\mathbf{e}_i, \mathbf{n}_i\})$  and  $f_7(\mathbf{T}_{id}^*, \{\mathbf{e}_i, \mathbf{n}_i\})$  between  $\{\mathbf{e}_i, \mathbf{n}_i\}$  and other variables are already in model 3.6, the new components are potentially three-way interactions.

My initial guess was to add a three-margin tensor product spline between  $\mathbf{T}_{id}^0$ ,  $\mathbf{T}_{id}^*$  and  $\{\mathbf{e}_i, \mathbf{n}_i\}$ , but the resulting model does not remove spatial autocorrelation in residuals. I then turned to the following model

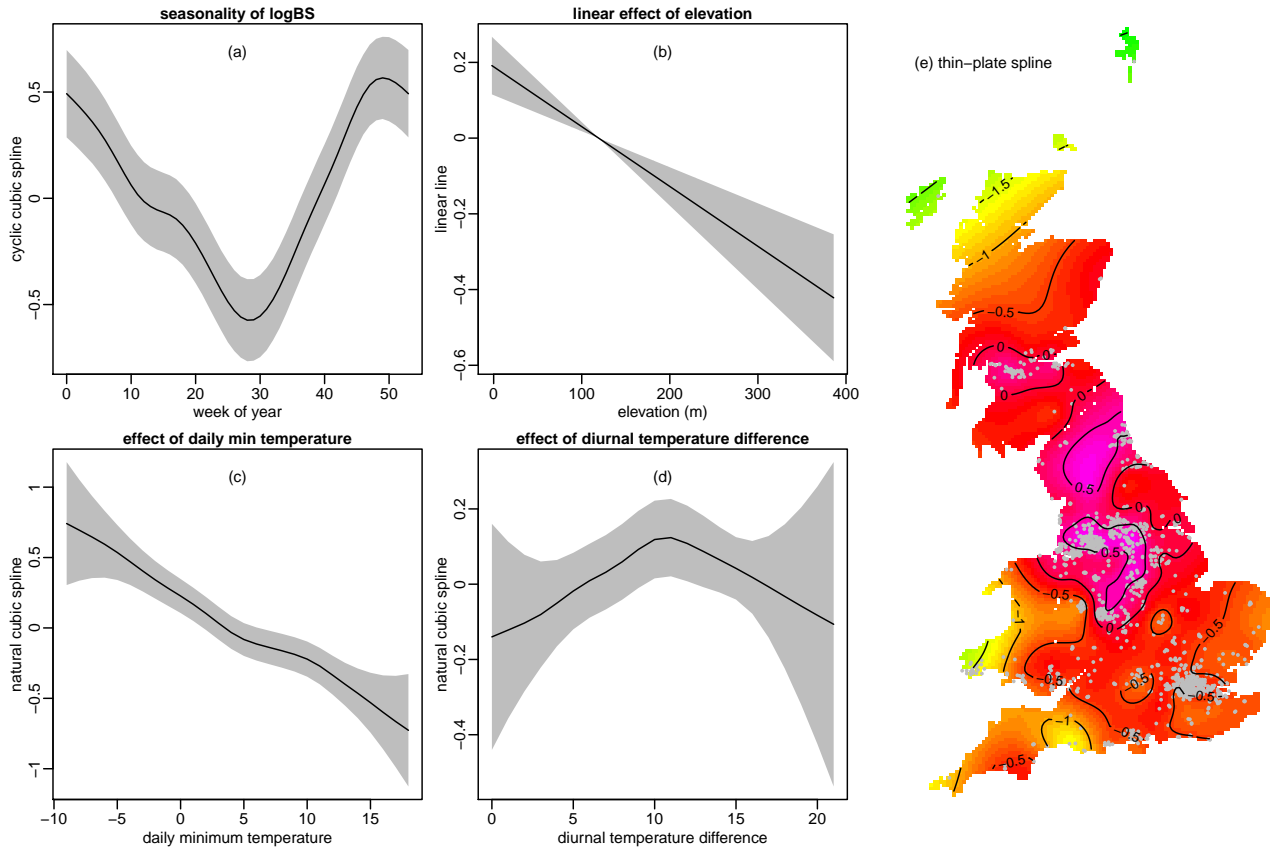


**Figure 3.41:** Validating i.i.d. assumption for weekly random effects in model 3.6. The left panel is a histogram of predicted random effects and the right panel is their sample autocorrelation. No temporal autocorrelation is spotted.

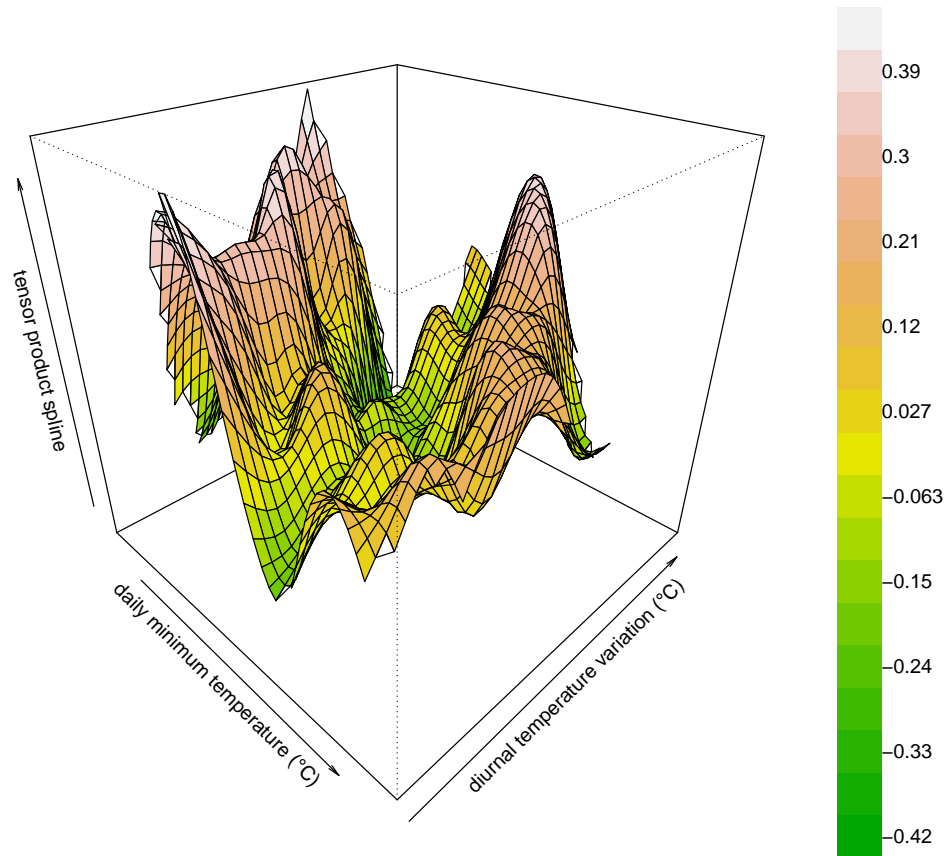


**Figure 3.42:** Validating i.i.d. assumption for station-specific random effects in model 3.6. The left panel is a histogram of predicted random effects and the right panel is their empirical variogram. No spatial autocorrelation is spotted.

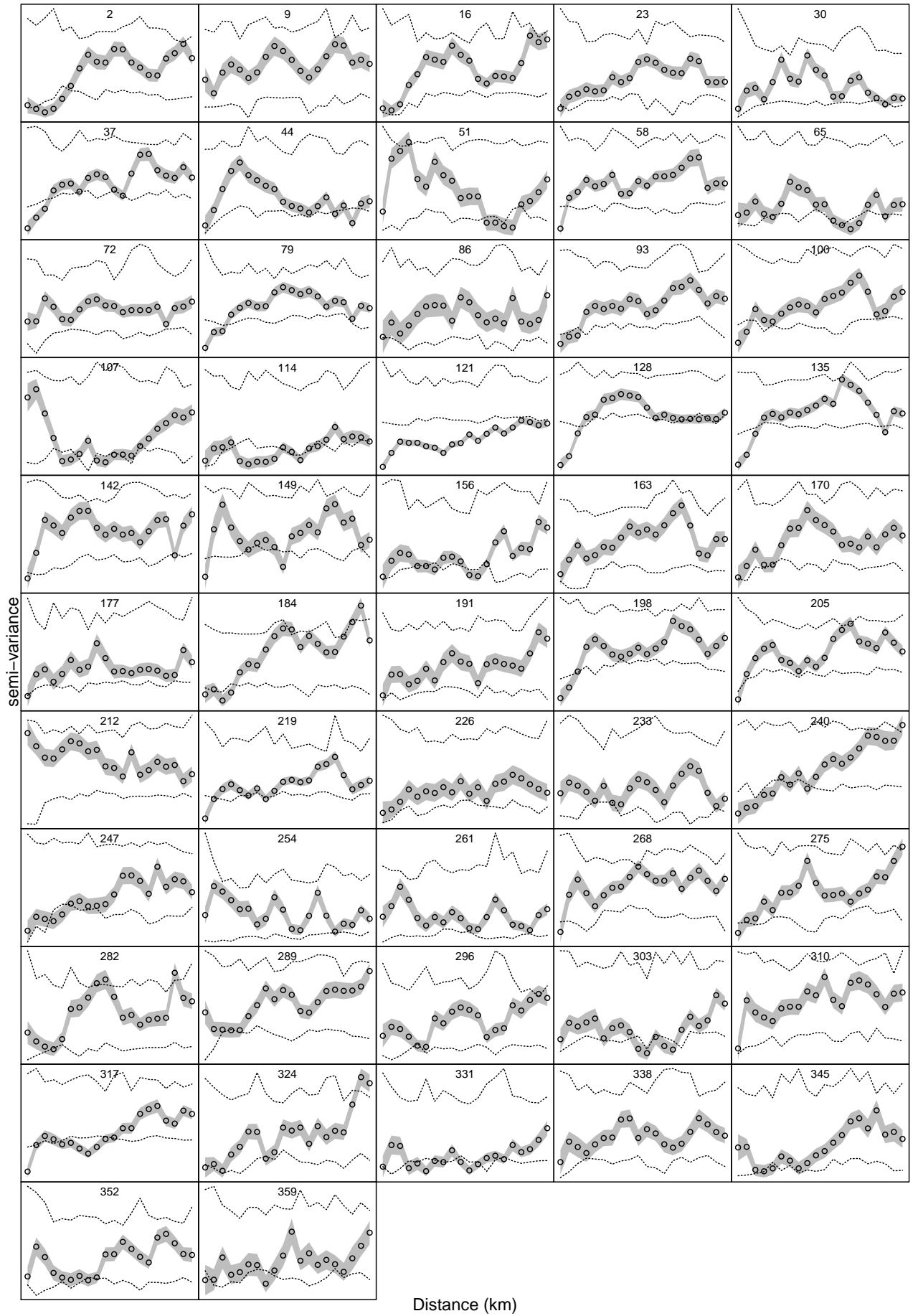




**Figure 3.43:** Estimated main effects in fitted model 3.6, namely  $\hat{f}_1$  (panel (a)),  $\hat{f}_2$  (panel (e)),  $\hat{f}_{11}$  (panel (b)),  $\hat{f}_4$  (panel (c)) and  $\hat{f}_6$  (panel (d)).



**Figure 3.44:** A perspective plot to visualize the tensor product spline  $\hat{f}_8$  in fitted model



**Figure 3.45:** Empirical variograms of spatial residuals from fitted model 3.6 on each of the 52 Mondays in year 1967. Axes are omitted for compact display. The variogram is computed up to a distance of 200km at every 5km bin starting from 0 (which is consistent with Figure 3.31). The day of year (1 ~ 365) is labelled on each plot.

**Table 3.10:** Summary of fitted model 3.7 to each day of week in 1967.

(a)  $n$  is the number of data;  $p$  is the number of regression coefficients;  $edf$  is the effective degree of freedom of the fitted model;  $\hat{\sigma}_{\text{iw}}^2$  is the estimated residual variance;  $\hat{\sigma}_9^2$  is the estimated variance for station-specific i.i.d. random effect  $\hat{f}_9$ ;  $\hat{\sigma}_{10}^2$  is the estimated variance for week of year i.i.d. random effect  $\hat{f}_{10}$ ; adj.  $R^2$  is the adjusted R-squared of the fitted model.

	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
$n$	54386	54581	55874	56612	56006	54437	55357
$p$	10348	10348	10347	10347	10347	10346	10349
$edf$	3918.89	4090.49	4046.72	4132.20	4270.01	4045.32	3893.77
$\hat{\sigma}_{\text{iw}}^2$	0.1837	0.1962	0.1923	0.1999	0.1910	0.1968	0.1869
adj. $R^2$	85.47%	83.89%	84.05%	83.02%	82.40%	84.16%	85.43%
$\hat{\sigma}_9^2$	0.1142	0.1211	0.1142	0.1154	0.1102	0.1069	0.1203
$\hat{\sigma}_{10}^2$	0.0957	0.0747	0.1303	0.1403	0.1489	0.1408	0.1145

(b) Termwise number of regression coefficients  $p$  and effective degree of freedom.

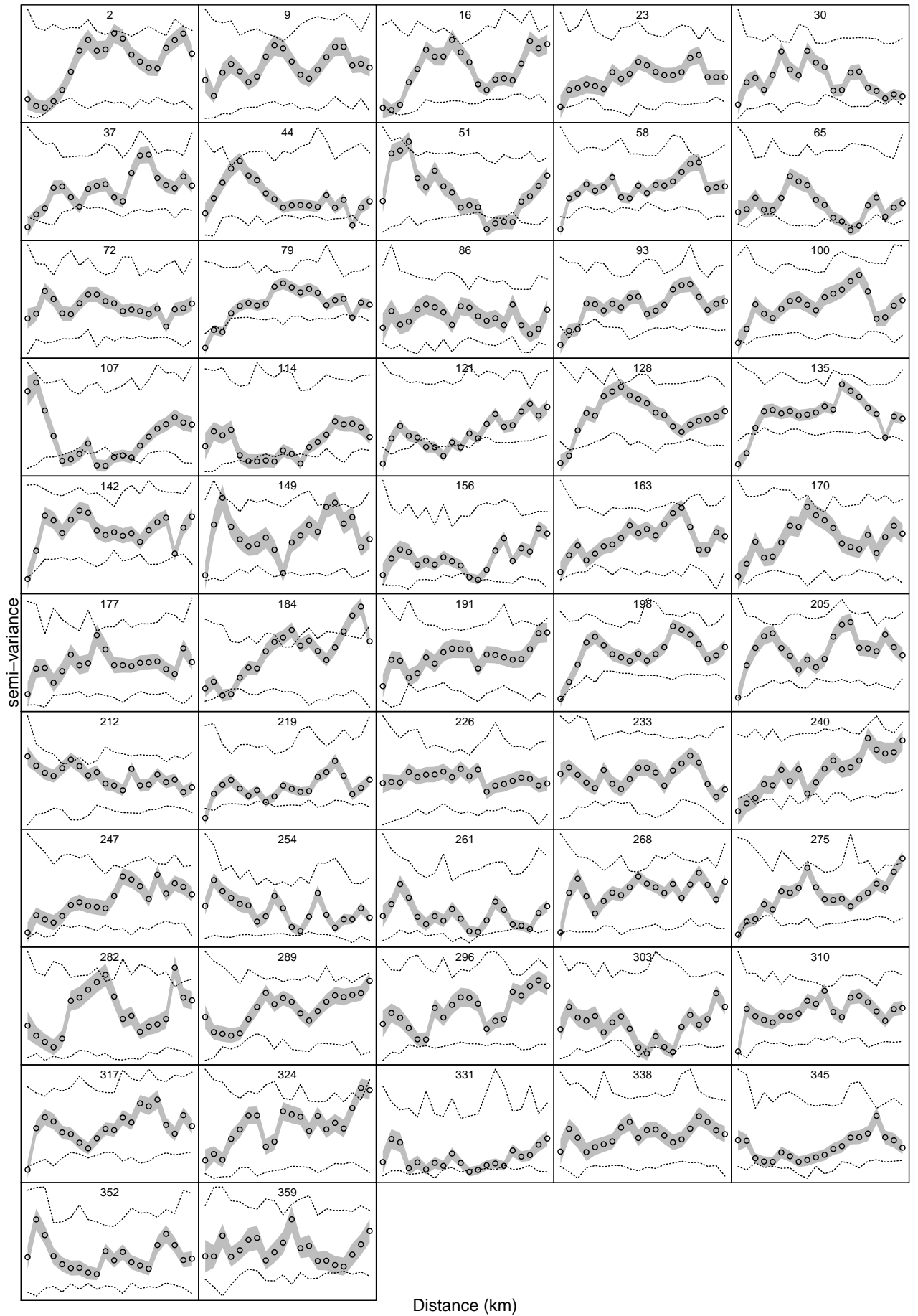
	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
$\hat{f}_0$ ( $p = 5$ )	3.98	3.97	3.97	3.97	3.97	3.97	3.98
$\hat{f}_1$ ( $p = 8$ )	4.28	4.35	4.03	2.87	2.56	3.31	3.28
$\hat{f}_2$ ( $p = 149$ )	95.46	94.91	98.26	93.93	103.10	100.06	93.56
$\hat{f}_3$ ( $p = 2682$ )	1294.29	1328.91	1344.74	1287.46	1461.13	1355.35	1206.16
$\hat{f}_4$ ( $p = 9$ )	4.17	5.69	2.73	5.38	5.80	5.24	3.28
$\hat{f}_5$ ( $p = 891$ )	196.16	202.92	255.95	230.91	286.09	246.49	212.08
$\hat{f}_6$ ( $p = 9$ )	1.44	1.24	1.00	1.00	1.00	4.00	4.23
$\hat{f}_7$ ( $p = 891$ )	210.96	243.02	173.71	260.18	202.78	221.31	232.16
$\hat{f}_8$ ( $p = 196$ )	56.01	61.43	57.07	60.98	39.05	35.45	49.98
$\hat{f}_9$ ( $p = 1275$ )	1119.04	1119.58	1115.44	1117.59	1109.85	1106.46	1123.40
$\hat{f}_{10}$ ( $p = 52$ )	46.25	46.05	46.63	47.75	48.07	47.33	48.35
$\hat{f}_{11}$ ( $p = 4$ )	1.00	1.00	1.00	1.00	1.00	1.00	1.00
$\hat{f}_{12}$ ( $p = 2088$ )	421.92	503.85	468.50	531.82	490.31	441.06	442.65
$\hat{f}_{13}$ ( $p = 2088$ )	462.90	472.57	472.69	486.35	514.29	473.28	468.66

**Model 3.7:** model 3.6 + three-way interaction

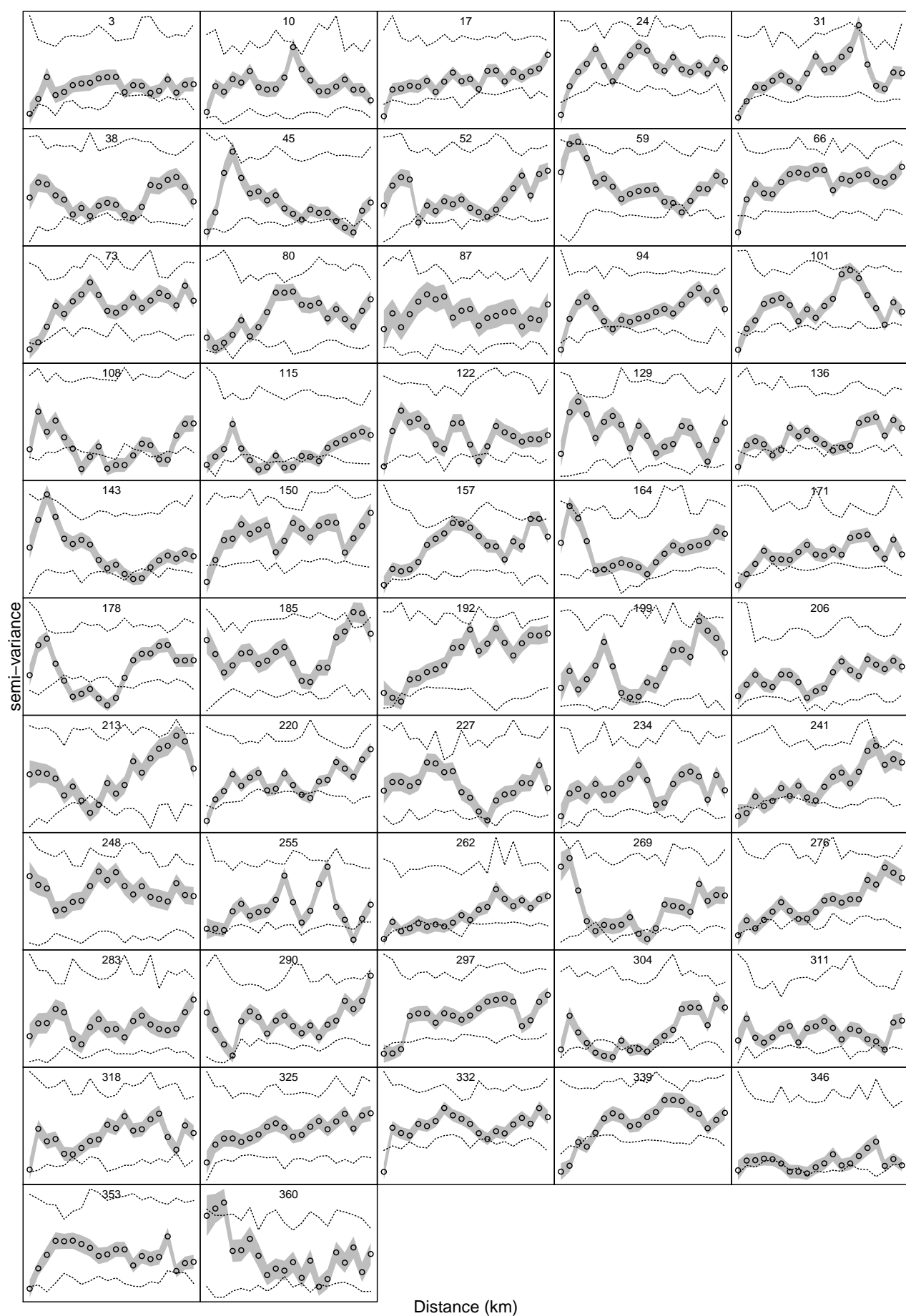
$$\begin{aligned}
\log\text{BS}_{\text{iw}} = & f_0(\mathbf{E}_i) + f_1(\mathbf{w}_y; 10) + f_2(\{\mathbf{e}_i, \mathbf{n}_i\}; 150) + f_3(\mathbf{w}_y, \{\mathbf{e}_i, \mathbf{n}_i\}; 20, 150) + \\
& f_4(\mathbf{T}_{\text{id}}^0; 10) + f_5(\mathbf{T}_{\text{id}}^0, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 100) + \\
& f_6(\mathbf{T}_{\text{id}}^*; 10) + f_7(\mathbf{T}_{\text{id}}^*, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 100) + \\
& f_8(\mathbf{T}_{\text{id}}^0, \mathbf{T}_{\text{id}}^*; 15, 15) + f_9(i; 1275) + f_{10}(\mathbf{w}_y; 52) + f_{11}(\mathbf{h}_i; 5) + \\
& f_{12}(\mathbf{w}_y, \mathbf{T}_{\text{id}}^0, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 10, 30) + f_{13}(\mathbf{w}_y, \mathbf{T}_{\text{id}}^*, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 10, 30) + \epsilon_{\text{iw}},
\end{aligned}$$

where two three-way interactions  $f_{12}$  and  $f_{13}$  between space, time and temperature are added. Note that compared with model 3.6, the  $k$  values of many terms have been reduced. In particular, the  $k$  value for  $\{\mathbf{e}_i, \mathbf{n}_i\}$  margin in  $f_{12}$  and  $f_{13}$  is set very low. Whether this is reasonable or not, it is the only way to be able to practically fit this model. Model 3.6 readily has more than 10000 coefficients, and if  $k$  values are not reduced, this new model would involve 110000+ coefficients, even more than the number of data (54386). After reducing  $k$ , this model still has about 10000 parameters. This model turns out to successfully remove all spatial autocorrelation in residuals (see Figure 3.46). I have also fitted this model to logBS from other days of week, and without exception, no spatial autocorrelation is seen from spatial residuals (see Figures 3.47, 3.48, 3.49, 3.50, 3.51, 3.52). Table 3.10 is a summary of those seven fitted models.

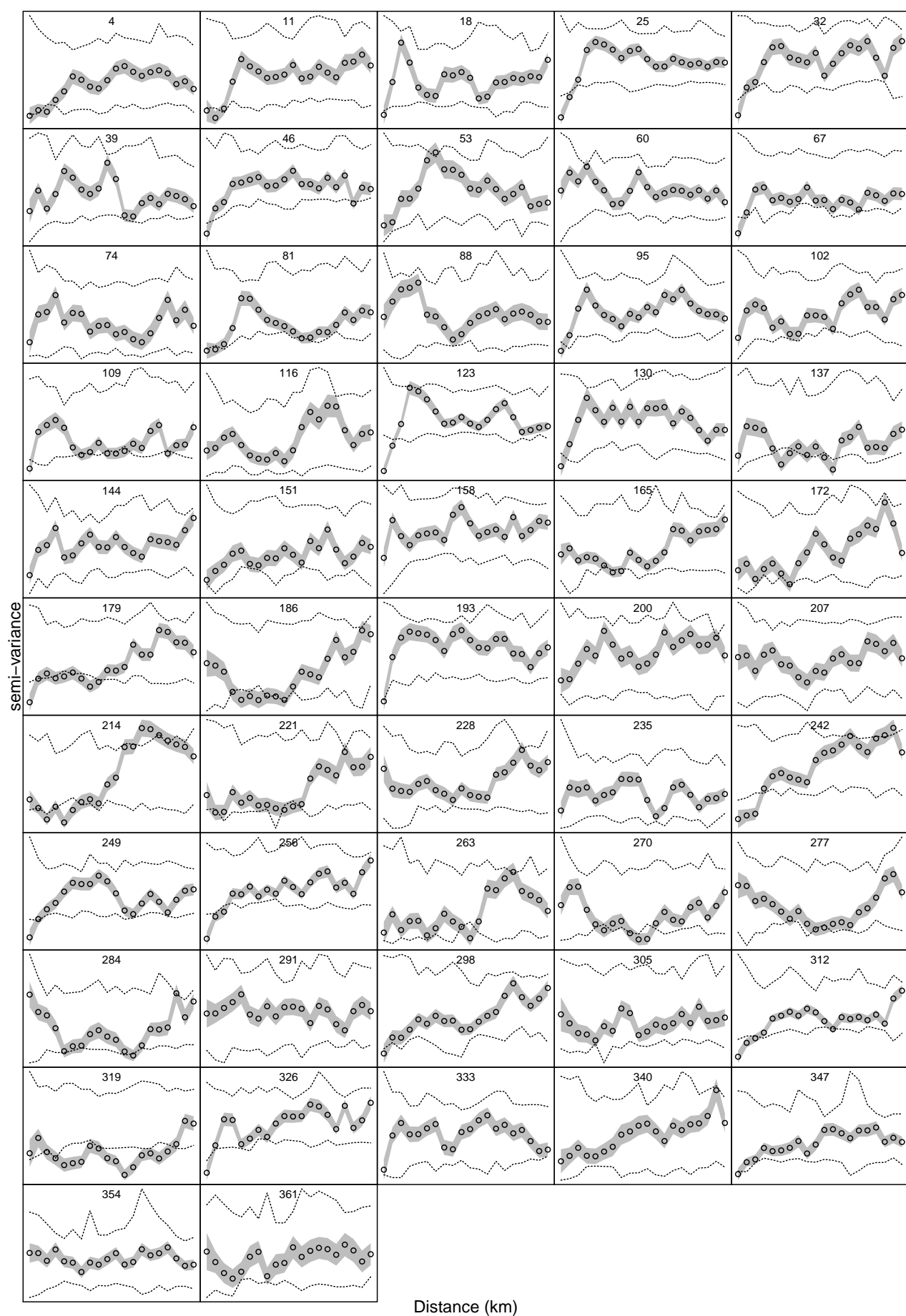
However, it is hard to interpret these three-way interaction. Basically they mean that temperature effects can change at any location and any time (hence essentially saying that there is no relationship between temperature and logBS?). In fact, if you think about this carefully, it is not difficult to understand why these components are so powerful to remove spatial autocorrelation. Temperature



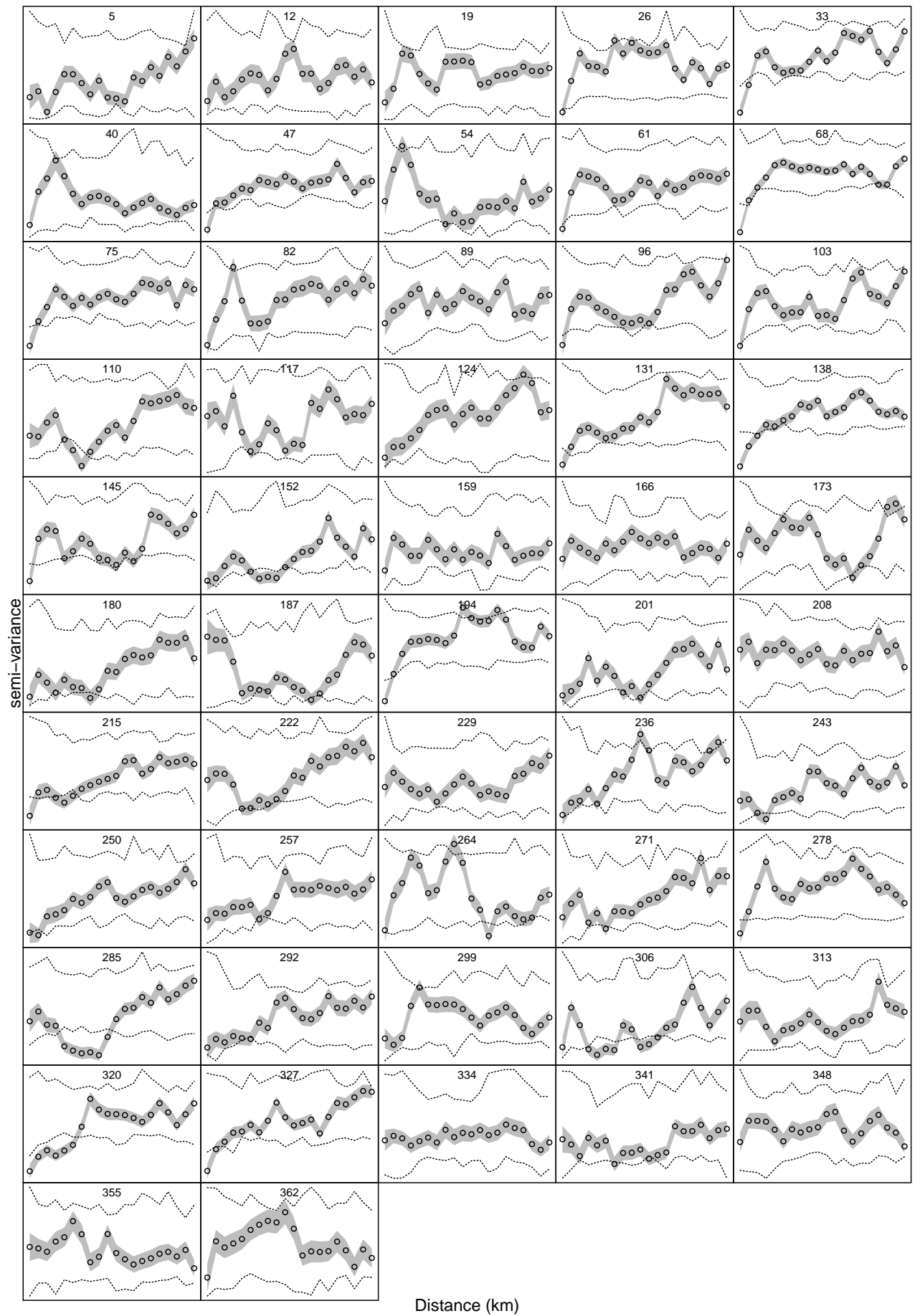
**Figure 3.46:** Empirical variograms of spatial residuals from fitted model 3.7 on each of the 52 **Mondays** in year 1967. Compared with Figure 3.46, the three-way interactions have helped remove all spatial autocorrelation in residuals.



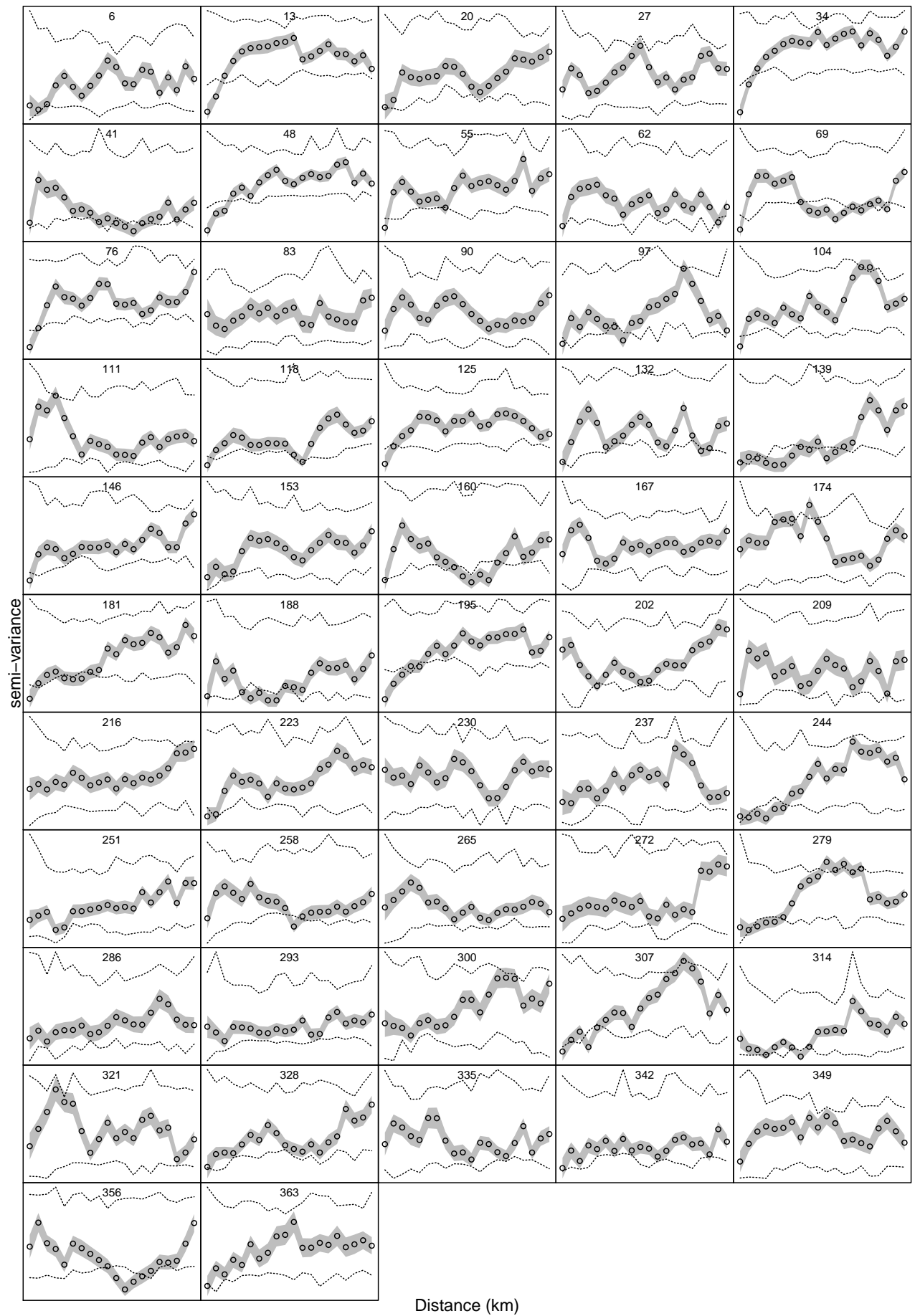
**Figure 3.47:** Empirical variograms of spatial residuals from fitted model 3.7 on each of the 52 **Tuesdays** in year 1967.



**Figure 3.48:** Empirical variograms of spatial residuals from fitted model 3.7 on each of the 52 **Wednesdays** in year 1967.

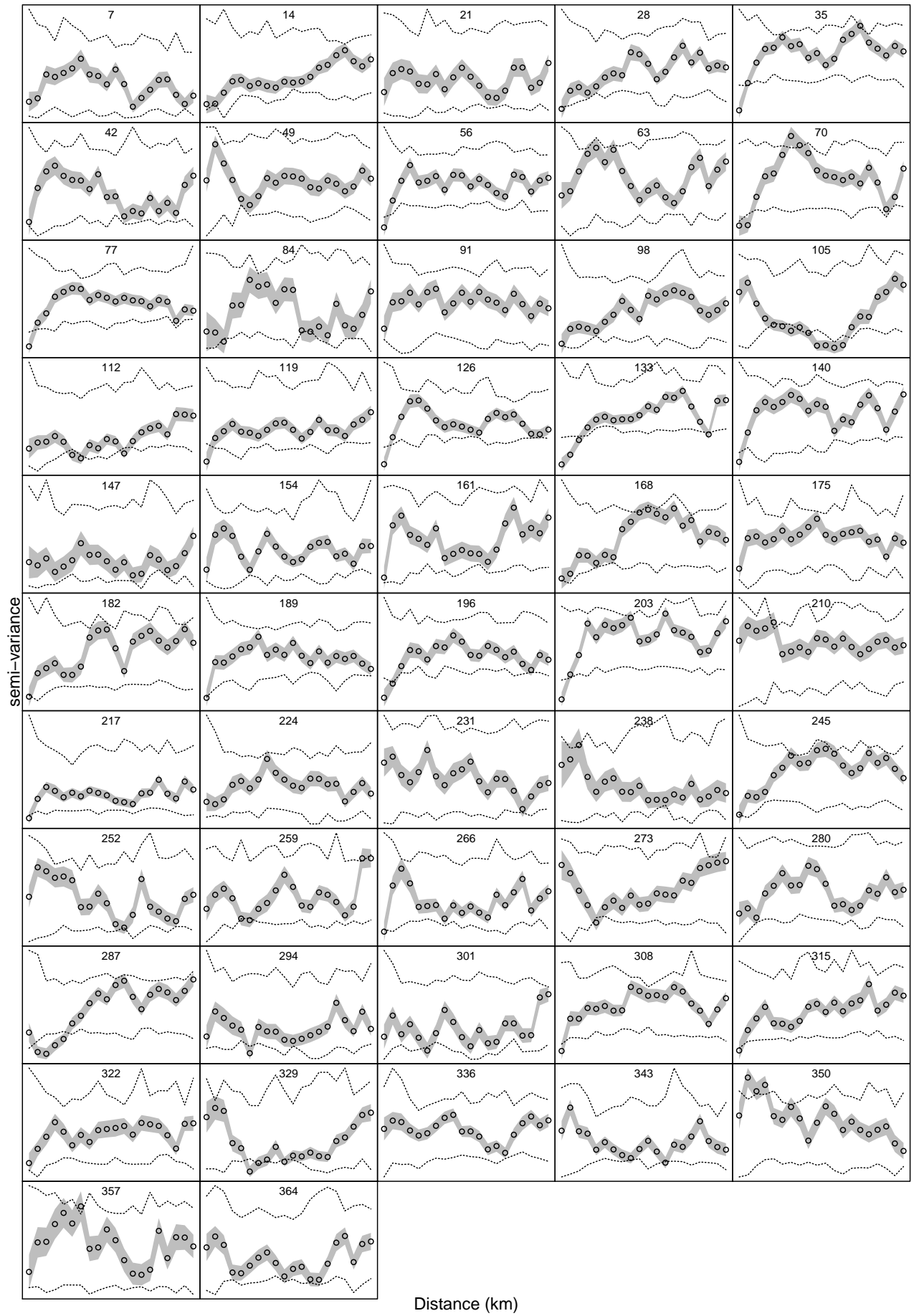


**Figure 3.49:** Empirical variograms of spatial residuals from fitted model 3.7 on each of the 52 **Thursday** in year 1967.

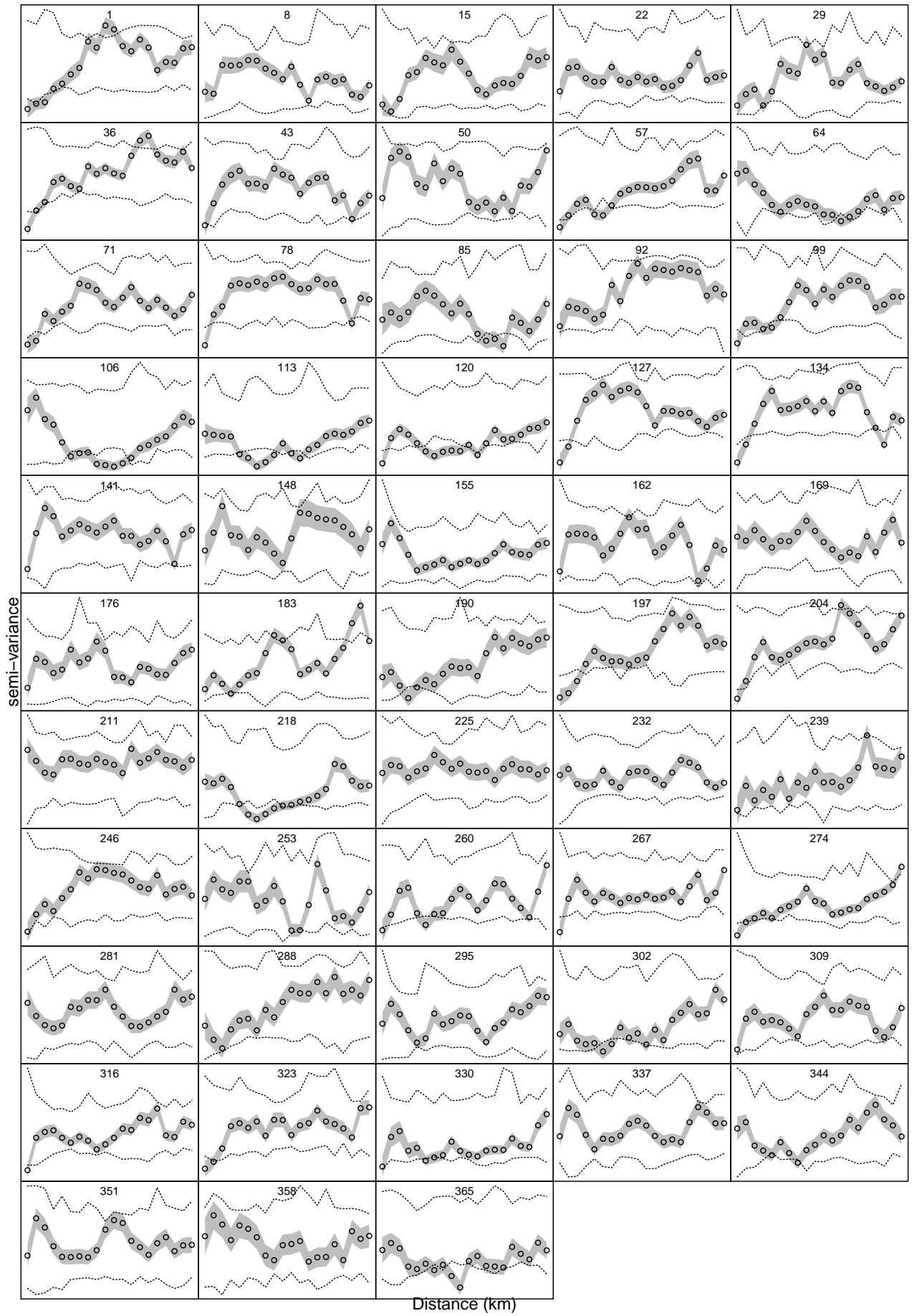


**Figure 3.50:** Empirical variograms of spatial residuals from fitted model 3.7 on each of the 52 **Friday** in year 1967.





**Figure 3.51:** Empirical variograms of spatial residuals from fitted model 3.7 on each of the 52 Saturdays in year 1967.



**Figure 3.52:** Empirical variograms of spatial residuals from fitted model 3.7 on each of the 53 **Sundays** in year 1967.

**Table 3.11:** Number of data  $n$ , model complexity  $p$  (number of parameters) and model fitting time  $T$  (in seconds) for logBS models so far.

<i>Manchester 11</i> model 3.3b			annual mean logBS model 3.5			Monday logBS in 1967 model 3.6		
n	p	T	n	p	T	n	p	T
15707	732	442.15	24239	6750	9927.68	54386	10932	45962.61

variables are the only daily variables that are available in my Black Smoke dataset. Even a fake linear model like

$$\text{logBS}_{\text{id}} = \alpha_{\text{id}}^0 \mathbf{T}_{\text{id}}^0 + \alpha_{\text{id}}^* \mathbf{T}_{\text{id}}^* + \epsilon_{\text{id}}$$

would be able to well “predict” logBS if the coefficients  $\alpha_{\text{id}}^0$  and  $\alpha_{\text{id}}^*$  can vary freely over space and time. So I am very skeptical of model 3.7.

The need for these odd-looking three-way interactions is probably a hint that other daily covariates than temperatures are necessary for modelling daily logBS. However, this can not be possibly done in this PhD research (acquiring daily temperature from UKCP09 had already been a non-trivial task).

Should these three-way interactions be used for model development? A possible way to answer this is running cross-validation. However, this is far beyond the computational capability that current GAM fitting method could offer.

### 3.5.5 Encountering computational hurdle

Previous presentation on model development was made so smoothly, as if they had been easy to obtain. The true story is, model fitting was very slow. Practical model building was further slowed down by the non-negligible costs for choosing  $k$  and updating a present model with new components. Table 3.11 tells you how much time was spent for fitting model 3.3b for *Manchester 11*, model 3.5 for annual mean logBS and model 3.6 for Monday logBS in 1967. The number of data and number of parameters / coefficients are also provided for a reference. Fitting time for model 3.7 is not included. It has similar  $n$  and  $p$  to model 3.6 hence similar fitting time. Note that fitting these daily models takes half a day each. It had taken nearly a week to produce the results in Table 3.10 and Figures 3.46 to 3.52. So it is time to do something to speed up model fitting, otherwise no further models can be feasibly built for Black Smoke.

## 3.6 Summary

The model development for logBS in this Chapter, on one hand, offers solid case studies for time series modelling and spatial-temporal modelling via GAM, and on the other hand, reveals the high complexity required for building GAMs for daily logBS even just from a single year. Many statistical questions can be raised regarding the model development, like

1. What is the reason behind the spatial-temporal extrapolation problem in model 3.5 and is there any way to suppress it?
2. Should three-way interactions in model 3.7 be used to adequately model space-time relationship? Won't they actually lead to overfitting?

But they all seem to give way to the computational hurdle in fitting “big” models: models with about  $10^5$  data and parameters.

It is time to dive into the computational engine of additive models and improve its design so that model estimation with big datasets and high complexity is feasible and more importantly, as efficient as possible.

To see opportunities in upgrading GAM computational engine requires in the first place a solid understanding of the existing one. So the next Chapter will review all computational details for GAM with large datasets. It is also at this point that the flavour of this thesis is going to change. There will be fewer nice looking pictures; instead, there will be more complicated mathematical notations, algorithms and benchmarking.

## Chapter 4

# GAM computation for large datasets: a review

This Chapter formally introduces the computational details underneath the `bam` function of R package `mgcv` (prior to `mgcv_1.8-7`), which is designed for estimating a GAM with large datasets.

GAM computations are ultimately about matrix computations, so from now on, issues of matrix computations are recurring. To ease later demonstrations, §4.1 will offer some preliminaries on the most frequently referenced matrix computations in this thesis.

Methods for large datasets was first developed for additive models. It breaks down the model fitting process into two stages: model matrix reduction (to be introduced in §4.2) and REML estimation (to be introduced in §4.3). Methods for generalized additive models is an iterative application of such method for additive models. Such way of estimating a GAM is termed “performance iteration”. There is another way for GAM estimation called “outer iteration”. It uses a different logic on iterations and is implemented in the `gam` function of `mgcv`. I will in §4.4 give a brief introduction to this method, pointing out its key difference from “performance iteration”. Computational details will not be provided as the method is infeasible for large datasets.

These computational details in this Chapter are largely based on Wood (2011) and Wood et al. (2015), but integrated together for a coherent presentation on estimation of additive models. Materials in this Chapter should be a very good reference for anyone who wants to dive into `bam`. The two papers essentially use different estimation methods: the former is implemented in `gam` function using “outer iteration” and is not “bigdata” oriented. The latter one mainly demonstrates the model matrix reduction idea, and simply cites the former for derivative computation in REML estimation. However, while they do share something similar, the derivative computation is still substantially different.

I would also like to give a pre-notice, that my demonstration for REML estimation in §5 is not 100% “loyal” to current `bam` implementation. In the learning process of the computational engine, I discovered a few places where the design idea could be better implemented. These actually contribute to my first round of performance optimization for `bam`, to be benchmarked in the next Chapter. In this Chapter, I will just present my improved version of REML estimation; the inefficiency in the original version will be pointed out in the next Chapter.

## 4.1 Preliminaries on matrix computations

This section gives a brief introduction to some important matrix computations in the heart of scientific computing.

- §4.1.1 introduces Householder transformation, Householder QR factorization, the mathematical meaning of QR factorization, the “factored form storage” of factor ‘Q’, and matrix multiplication involving factor ‘Q’.
- §4.1.2 introduces Cholesky factorization and its link with QR factorization.
- §4.1.3 and §4.1.4 introduce some computations associated with triangular matrices, namely solving a triangular system of linear equations, inverting a triangular matrix and computing cross-product of a triangular matrix. These operations are essential since both QR factorization and Cholesky factorization produce a triangular matrix factor.
- §4.1.5 introduces a subtype of general eigen decomposition: the symmetric eigen decomposition for a real symmetric matrix.

I will not demonstrate algorithms for all computations. The aim is to provide background information for their practical applications in GAM computations to assist you in understanding subsequent sections in this Chapter. Basic knowledge of linear algebra like matrix multiplication, matrix inverse, positive-definite matrix and orthogonal matrix (rotation matrix or reflection matrix) is a prerequisite.

Numerical matrix computations primarily consist of floating-point arithmetic, specially floating-point multiplication and addition, hence their computational complexity is often measured by the amount of multiplication and addition they involve, known as the number of floating-point operations (FLOP). For example, matrix multiplication between an  $m \times k$  matrix  $\mathbf{A}$  and a  $k \times n$  matrix  $\mathbf{B}$  involves  $2mkn$  FLOP, since 1) the resulting matrix has  $mn$  elements; 2) element  $(i, j)$  is computed by a dot product  $\sum_{l=1}^k \mathbf{A}(i, l)\mathbf{B}(l, j)$  which involves  $k$  addition and  $k$  multiplication. FLOP count for all matrix computations covered in this section will be given along the way, but is also summarized in §4.1.6. The summary section also gives numerous fruitful examples on how matrix factorizations are useful in practice. You may have a quick glimpse on that section first to gain extra motivation to go through the following possibly tedious pages.

### 4.1.1 Householder QR factorization

There are several types of QR factorization, but a particular type known as Householder QR factorization is applied throughout this thesis.

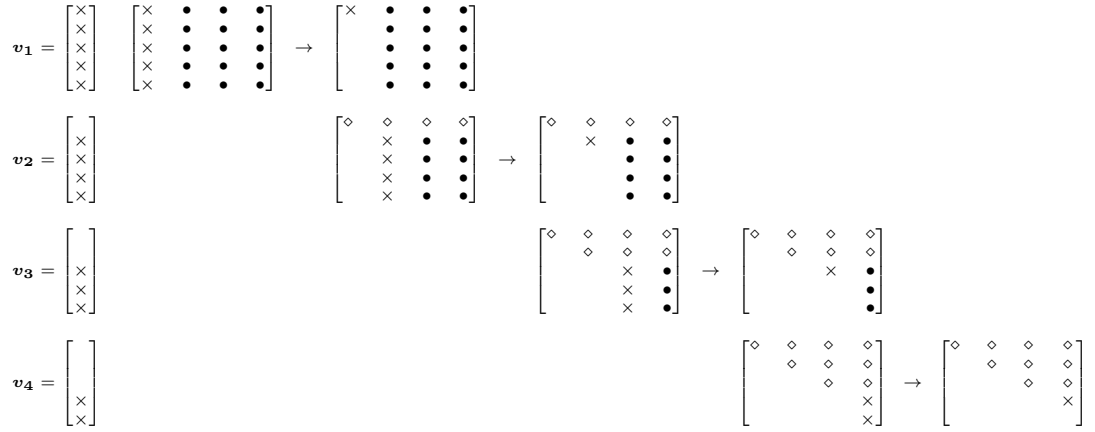
A Householder matrix is a reflection matrix  $\mathbf{H} = \mathbf{I} - \tau \mathbf{v} \mathbf{v}'$ , where  $\mathbf{v}$  is a non-zero vector called a reflector vector and  $\tau = \frac{2}{\|\mathbf{v}\|^2}$ . Given two length- $n$  vectors  $\mathbf{x}$  and  $\mathbf{y}$ , there exists a Householder transformation  $\mathbf{H}\mathbf{x} = \mathbf{y}$ , where  $\mathbf{v} = \mathbf{x} - \mathbf{y}$ . A special type of Householder transformation is

$$\mathbf{H} : \begin{bmatrix} \mathbf{x}(1) \\ \mathbf{x}(2:n) \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{y}(1) \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{y}(1) = -\text{sgn}(\mathbf{x}(1))\sqrt{\mathbf{x}(1)^2 + \|\mathbf{x}(2:n)\|^2}.$$

Such type of transformation can be sequentially applied to zero the lower triangular part of a matrix, yielding an upper triangular matrix. See Figure 4.1 for an illustration.

In general, let  $\mathbf{X}$  be an  $n \times p$  matrix ( $n \geq p$ ),  $p$  sequential Householder transformations to  $\mathbf{X}$  yields an upper triangular matrix

$$\mathbf{H}_p \mathbf{H}_{p-1} \cdots \mathbf{H}_1 \mathbf{X} = \mathbf{R}.$$



**Figure 4.1:** Illustration of Householder QR factorization using a  $5 \times 4$  matrix. Four Householder transformations are needed to arrive at a final upper triangular matrix. Zero entries in matrices and vectors are left blank. Place holders  $\bullet$ ,  $\times$  and  $\diamond$  represent non-zero entries, with the following meanings:  $\times$  are used to derive reflector vectors;  $\bullet$  are updated by the resulting Householder transformation;  $\diamond$  are unaffected as the “range” of the transformation is shrinking toward the lower right submatrix at each step.

These transformations can be summarized by a single orthonormal transformation  $\mathbf{Q}'\mathbf{X} = \mathbf{R}$ , or  $\mathbf{X} = \mathbf{Q}\mathbf{R}$ , which is a QR factorization. Note that

- The  $n \times n$  matrix  $\mathbf{Q} = \mathbf{H}_1\mathbf{H}_2 \cdots \mathbf{H}_p$  is mostly useful for paper work. It needs not be formed and is practically stored in a “factored form” as a number of reflector vectors. Later I will explain the storage of such “factored form” and how matrix multiplication involving  $\mathbf{Q}$  is practically done;
- The  $\mathbf{R}$  matrix above is not a triangular matrix, but its upper  $p \times p$  submatrix is. In many theoretical derivations, QR factorization is often written in a partitioned or “thin” form

$$\mathbf{X} = (\mathbf{Q} \quad \mathbf{Q}^\perp) \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} = \mathbf{Q}\mathbf{R},$$

where  $\mathbf{Q} = \mathbf{Q}(:, 1:p)$ ,  $\mathbf{Q}^\perp = \mathbf{Q}(:, (p+1):n)$  and  $\mathbf{R} = \mathbf{R}(1:p, :)$ .

For more details on Householder transformation and Householder QR factorization, see (Golub and Loan, 2013, §5.1 and §5.2).

If  $\mathbf{A}$  has full column rank, then finding an orthonormal basis for the column space of  $\mathbf{A}$  is a Gram-Schmidt process. In fact, QR factorization is the matrix representation for Gram-Schmidt process (just like LU factorization is that for Gaussian elimination), and the columns of  $\mathbf{Q}$  is just the desired basis. At the same time, columns of  $\mathbf{Q}^\perp$  is an orthonormal basis for the null space (or kernel) of  $\mathbf{A}'$  (A more concrete explanation will be made in §4.1.6).

But such nice interpretation of QR factorization does not hold when  $\mathbf{A}$  does not have full column rank. A column pivoting strategy is used to overcome this and support rank estimation. The resulting thin-QR factorization has the form

$$\mathbf{A}\mathbf{P} = \mathbf{Q}\mathbf{R},$$

where  $\mathbf{P}$  is a column permutation matrix (a variant of an identity matrix by column permutation), and  $\mathbf{R}$  is an upper triangular matrix with  $|\mathbf{R}(1, 1)| \geq |\mathbf{R}(2, 2)| \geq \cdots \geq |\mathbf{R}(p, p)| \geq 0$ . In practice,  $\mathbf{P}$  is simply stored as an index vector  $\mathbf{k}$ , so that  $(\mathbf{A}\mathbf{P})(i) = \mathbf{A}(:, \mathbf{k}(i))$ . Golub and Loan (2013, §5.4.2) gives algorithm details for such factorization. Note that the factorization helps supporting rank estimation but is not directly rank-revealing. While it is tempting to use  $|\mathbf{R}(i, i)| / |\mathbf{R}(1, 1)| < \epsilon$  for rank estimation, where  $\epsilon$  is some numerical tolerance, Golub and Loan (2013, §5.4.3) gives a counterexample. The condition estimation described in Golub and Loan (2013, §3.5.4) can serve for

rank estimation. If  $\mathbf{R}$  has an estimated rank of  $r$ , then columns of  $\mathbf{Q}(1:r)$  is an orthonormal basis of the column space of  $\mathbf{A}$ .

QR factorization with column pivoting (hereafter called “pivoted QR factorization” for short) can be seen as a QR factorization with a non-triangular ‘R’ factor. A column permutation matrix is invertible, and in fact is orthogonal with  $\mathbf{P}^{-1} = \mathbf{P}'$ . Therefore,

$$\mathbf{AP} = \mathbf{QR} \Rightarrow \mathbf{A} = \mathbf{QR}^*, \mathbf{R}^* = \mathbf{RP}'.$$

Applying  $\mathbf{P}'$  to a matrix from the right is reversing the column permutation:  $\mathbf{R}^*(k(i)) = \mathbf{R}(i)$ . Later in this thesis, a thin QR factorization may be expressed by  $\mathbf{A} = \mathbf{QR}^*$  for brevity, and you should remember that  $\mathbf{R}^*$  in this case absorbs the reverse column permutation.

QR factorization has a computational complexity of  $(2np^2 - \frac{2}{3}p^3)$  FLOP. Some literature dealing with “tall-thin” matrices with  $n \gg p$ , often put it  $2np^2$ . However, if  $n = p$ , the complexity is  $\frac{4}{3}p^3$ . Pivoting would add some  $O(np)$  overhead, but it is ignorable.

You may have already spotted from Figure 4.1 that (non-zero entries of) the reflector vectors  $\mathbf{v}_i$  can be stored into the lower trapezoidal part of  $\mathbf{A}$ , so that in the end of the factorization, the upper triangular part of  $\mathbf{A}$  gives  $\mathbf{R}$ , while the lower trapezoidal part holds the sequence of  $\mathbf{v}_i$ . However, both  $\mathbf{v}_i$  and  $\mathbf{R}$  demand the diagonal of  $\mathbf{A}$ . Such clash can be avoided if  $\mathbf{v}_i$  is rescaled so that  $\mathbf{v}_i(i) = 1$ , then there is no need to store those 1. Such storage is called the “factored form storage” of  $\mathbf{Q}$ .

Matrix multiplications involving  $\mathbf{Q}$ , namely  $\mathbf{Q}'\mathbf{B}$ ,  $\mathbf{QB}$ ,  $\mathbf{BQ}'$  and  $\mathbf{BQ}$ , should be done by sequential Householder transformations. As  $\mathbf{BQ}' = (\mathbf{QB})'$  and  $\mathbf{BQ} = (\mathbf{Q}'\mathbf{B})'$ , it is sufficient to demonstrate the computation of left-side orthonormal transformations  $\mathbf{Q}'\mathbf{B}$  and  $\mathbf{QB}$ . For example, orthonormal transformation  $\mathbf{Q}'\mathbf{B}$  is just  $\mathbf{H}_p\mathbf{H}_{p-1}\cdots\mathbf{H}_1\mathbf{B}$ . Starting at  $\mathbf{B}_0 = \mathbf{B}$ , we iteratively compute  $\mathbf{B}_i = \mathbf{H}_i\mathbf{B}_{i-1}$ , and in the end  $\mathbf{Q}'\mathbf{B} = \mathbf{B}_p$ . Also note that there is even no need to form  $\mathbf{H}_i$ , as  $\mathbf{H}_i\mathbf{B} = (\mathbf{I} - \tau_i\mathbf{v}_i\mathbf{v}_i')\mathbf{B} = \mathbf{B} - \tau_i\mathbf{v}_i(\mathbf{B}'\mathbf{v}_i)'$  so the computation is not about matrix-matrix multiplication at all. Finally, the structure that  $\mathbf{v}_i(1:(i-1)) = \mathbf{0}$  can be exploited, so that  $\mathbf{H}_i\mathbf{B}$  would leave  $\mathbf{B}(1:(i-1), )$  unchanged. This explains why the “range” of Householder transformation shrinks in Figure 4.1 as the factorization proceeds.

Matrix multiplication involving  $\mathbf{Q}$  or  $\mathbf{Q}^\perp$  needs a post-processing or pre-processing, as

$$\begin{aligned}\mathbf{Q}'\mathbf{B} &= [\mathbf{Q} \quad \mathbf{Q}^\perp]' \mathbf{B} = \begin{bmatrix} \mathbf{Q}'\mathbf{B} \\ \mathbf{Q}^{\perp'}\mathbf{B} \end{bmatrix}, \\ \mathbf{QB} &= [\mathbf{Q} \quad \mathbf{Q}^\perp] \begin{bmatrix} \mathbf{B} \\ \mathbf{0} \end{bmatrix} = \mathbf{Q} \begin{bmatrix} \mathbf{B} \\ \mathbf{0} \end{bmatrix}, \\ \mathbf{Q}^\perp\mathbf{B} &= [\mathbf{Q} \quad \mathbf{Q}^\perp] \begin{bmatrix} \mathbf{0} \\ \mathbf{B} \end{bmatrix} = \mathbf{Q}^\perp \begin{bmatrix} \mathbf{0} \\ \mathbf{B} \end{bmatrix}.\end{aligned}$$

Thus for example,  $\mathbf{Q}'\mathbf{B} = (\mathbf{Q}'\mathbf{B})(1:p, )$  and  $\mathbf{QB}$  is an orthonormal transformation on a zero padded  $\mathbf{B}$ .

#### 4.1.2 Cholesky factorization

Cholesky factorization is a special case of LU factorization, which is a triangular matrix factorization for a full-rank square matrix. Let  $\mathbf{S}$  be a  $p \times p$  positive-definite matrix, its Cholesky factorization is

$$\mathbf{S} = \mathbf{LL}' \quad \text{or} \quad \mathbf{S} = \mathbf{R}'\mathbf{R},$$

where  $\mathbf{L}$  and  $\mathbf{R}$  are a lower and upper triangular matrix respectively, often called the lower and upper Cholesky factor of  $\mathbf{S}$ . An understanding of computation algorithm is useful but not essential to follow



this thesis. I refer keen readers to Golub and Loan (2013, §4.2.5) for the algorithm which involves  $\frac{1}{3}p^3$  FLOP.

The representation with upper triangular factor is most illuminating for establishing a link between Cholesky factorization and QR factorization, because if  $\mathbf{A} = \mathbf{QR}$ , then  $\mathbf{A}'\mathbf{A} = \mathbf{R}'\mathbf{R}$  so  $\mathbf{R}$  is “almost” the upper Cholesky factor of  $\mathbf{A}'\mathbf{A}$ . I used “almost” as in practice they are only identical up to sign flipping. Cholesky factorization algorithm yields an  $\mathbf{R}$  matrix with positive diagonal elements, while QR factorization does not, though it can be required (Demmel et al., 2009).

If  $\mathbf{S}$  is only positive-semidefinite, then the Cholesky factorization can only proceed with pivoting. This analogizes the QR factorization where  $\mathbf{A}$  does not have full column rank, because then  $\mathbf{A}'\mathbf{A}$  will be positive-semidefinite. The factorization in this case (hereafter called “pivoted Cholesky factorization”) has the form

$$\mathbf{P}'\mathbf{S}\mathbf{P} = \mathbf{L}\mathbf{L}' \quad \text{or} \quad \mathbf{P}'\mathbf{S}\mathbf{P} = \mathbf{R}'\mathbf{R},$$

where  $\mathbf{P}$  is a column permutation matrix (applying  $\mathbf{P}'$  to a matrix from the left is effectively doing row permutation), which is practically stored as an index vector, and  $\mathbf{R}$  is an upper triangular matrix with  $\mathbf{R}(1, 1) \geq \mathbf{R}(2, 2) \geq \dots \geq \mathbf{R}(p, p) \geq 0$ .

$\mathbf{S}$  needs better be pre-conditioned for greater numerical stability in the pivoted Cholesky factorization. This is particularly true when we factorize  $\mathbf{A}'\mathbf{A}$  as an analogy to QR factorization of  $\mathbf{A}$ , as  $\mathbf{A}'\mathbf{A}$  has twice as big the conditional number as  $\mathbf{A}$ . A common and simple pre-conditioner is the Jacobian diagonal pre-conditioner. For a symmetric matrix the pre-conditioner is applied from both sides of the matrix. Let  $\mathbf{J} = \text{diag}(\sqrt{\text{diag}(\mathbf{S})})$ , the pre-conditioned matrix is  $\mathbf{J}^{-1}\mathbf{S}\mathbf{J}^{-1}$ . The pivoted Cholesky factorization, for example, is then

$$\mathbf{P}'(\mathbf{J}^{-1}\mathbf{S}\mathbf{J}^{-1})\mathbf{P} = \mathbf{R}'\mathbf{R}.$$

Pivoted Cholesky factorization (and optionally with pre-conditioning) can be seen as a Cholesky factorization with non-triangular Cholesky factor, because

$$\begin{aligned} \mathbf{P}'\mathbf{S}\mathbf{P} = \mathbf{R}'\mathbf{R} &\Rightarrow \mathbf{S} = \mathbf{R}^*\mathbf{R}^*, \mathbf{R}^* = \mathbf{R}\mathbf{P}', \\ \mathbf{P}'(\mathbf{J}^{-1}\mathbf{S}\mathbf{J}^{-1})\mathbf{P} = \mathbf{R}'\mathbf{R} &\Rightarrow \mathbf{S} = \mathbf{R}^*\mathbf{R}^*, \mathbf{R}^* = \mathbf{R}\mathbf{P}'\mathbf{J}. \end{aligned}$$

Later in this thesis, a Cholesky factorization may be expressed by  $\mathbf{S} = \mathbf{R}^*\mathbf{R}^*$  for brevity, and you should remember that  $\mathbf{R}^*$  in this case absorbs the reverse column permutation and pre-conditioning.

### 4.1.3 Solving a triangular system of linear equations

Taking an upper triangular system  $\mathbf{R}\mathbf{x} = \mathbf{y}$ , or

$$\begin{pmatrix} R_{11} & R_{12} & \cdots & R_{1p} \\ & R_{22} & \cdots & R_{2p} \\ & & \ddots & \vdots \\ & & & R_{pp} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{pmatrix}$$

as an example, the solution is a back substitution

$$x_i = R_{ii}^{-1} \left( y_i - \sum_{k=i+1}^p R_{ik}x_k \right), \quad i = p, (p-1), \dots, 1.$$

That is,  $x_1, x_2, \dots, x_p$  can be determined one by one backward. For lower triangular system  $\mathbf{R}'\mathbf{x} = \mathbf{y}$  the solution is similarly a forward substitution. Both substitutions cost  $p^2$  FLOP, and they are practically used for computing  $\mathbf{x} = \mathbf{R}^{-1}\mathbf{y}$  or  $\mathbf{x} = \mathbf{R}'^{-1}\mathbf{y}$  without explicitly forming  $\mathbf{R}^{-1}$  or  $\mathbf{R}'^{-1}$ .

i = 0	i = 1	i = 2	i = 3
$\begin{pmatrix} \tilde{R}_{11} & \bullet & \bullet & \bullet \\ & \bullet & \bullet & \bullet \\ & & \bullet & \bullet \\ & & & \bullet \end{pmatrix}$	$\begin{pmatrix} \tilde{R}_{11} & \tilde{R}_{12} & \bullet & \bullet \\ & \tilde{R}_{22} & \bullet & \bullet \\ & & \bullet & \bullet \\ & & & \bullet \end{pmatrix}$	$\begin{pmatrix} \tilde{R}_{11} & \tilde{R}_{12} & \tilde{R}_{13} & \bullet \\ & \tilde{R}_{22} & \tilde{R}_{23} & \bullet \\ & & \tilde{R}_{33} & \bullet \\ & & & \bullet \end{pmatrix}$	$\begin{pmatrix} \tilde{R}_{11} & \tilde{R}_{12} & \tilde{R}_{13} & \tilde{R}_{14} \\ & \tilde{R}_{22} & \tilde{R}_{23} & \tilde{R}_{24} \\ & & \tilde{R}_{33} & \tilde{R}_{34} \\ & & & \tilde{R}_{44} \end{pmatrix}$

**Figure 4.2:** Illustration of triangular matrix inverse with overwriting using a  $4 \times 4$  matrix. In each iteration, elements in diagonal blocks and rectangular block are respectively coloured in black and gray. Those elements not used in the iteration is denoted by  $\bullet$ .

Vectors  $\mathbf{x}$  and  $\mathbf{y}$  can be generalized to  $p \times k$  matrices  $\mathbf{X}$  and  $\mathbf{Y}$ , then the problem is just  $k$  independent triangular system of equations, at a complexity of  $kp^2$  FLOP.

Later in this thesis, a pivoted triangular system is  $\mathbf{R}^* \mathbf{x} = \mathbf{y}$  or  $\mathbf{R}^{*'} \mathbf{x} = \mathbf{y}$  is often seen, where  $\mathbf{R}^*$  is the non-triangular factor from a QR or pivoted Cholesky factorization (and even pre-conditioning). However, this imposes no extra difficulty in equation solving. For example, in a pivoted Cholesky factorization with pre-conditioning, there is  $\mathbf{R}^* = \mathbf{R} \mathbf{P}' \mathbf{J}$ , where  $\mathbf{R}$  is truly triangular. Then

$$\begin{aligned} \mathbf{R}^* \mathbf{x} = \mathbf{y} &\Rightarrow \mathbf{x} = \mathbf{R}^{*-1} \mathbf{y} = \mathbf{J}^{-1} \mathbf{P} (\mathbf{R}^{-1} \mathbf{y}), \\ \mathbf{R}^{*'} \mathbf{x} = \mathbf{y} &\Rightarrow \mathbf{x} = \mathbf{R}^{*-1} \mathbf{y} = \mathbf{R}^{-1} (\mathbf{P}' \mathbf{J}^{-1} \mathbf{y}), \end{aligned}$$

so that the triangular structure of  $\mathbf{R}$  can be exploited for ordinary back and forward substitution.

#### 4.1.4 Positive-definite matrix inverse

Given a Cholesky factorization  $\mathbf{S} = \mathbf{R}' \mathbf{R}$  for a  $p \times p$  positive-definite matrix  $\mathbf{S}$ , the matrix inverse  $\mathbf{S}^{-1} = \mathbf{R}^{-1} \mathbf{R}'^{-1}$  can be obtained by first getting triangular matrix inverse  $\tilde{\mathbf{R}} = \mathbf{R}^{-1}$ , then doing a matrix cross-product  $\mathbf{S}^{-1} = \tilde{\mathbf{R}} \tilde{\mathbf{R}}'$ .

A direct way to compute  $\tilde{\mathbf{R}}$  is to solve the triangular system  $\mathbf{R} \tilde{\mathbf{R}} = \mathbf{I}$ . This appears to involve  $p^3$  FLOP but that  $\mathbf{I}$  is identity can make back substitution more efficient. The triangular system for  $\tilde{\mathbf{R}}(:, j)$  can be partitioned as

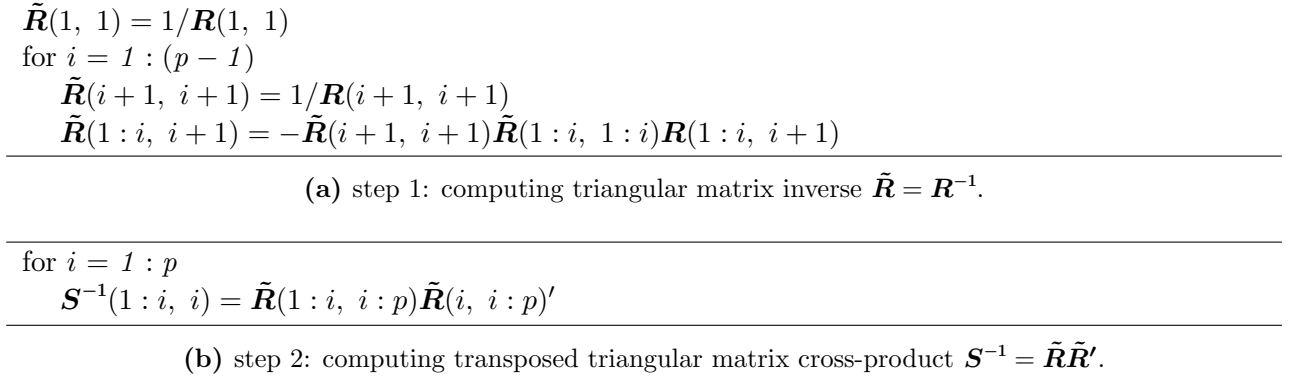
$$\begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ & \mathbf{R}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{z} \\ \mathbf{0} \end{pmatrix},$$

where  $\mathbf{z}(1 : (j-1)) = \mathbf{0}$  and  $\mathbf{z}(j) = 1$ . Since all diagonal elements of  $\mathbf{R}$  are positive, both  $\mathbf{R}_{11}$  and  $\mathbf{R}_{22}$  are invertible, thus  $\tilde{\mathbf{R}}((j+1) : p, j) = \mathbf{x}_2 = \mathbf{0}$  and  $\tilde{\mathbf{R}}(1 : j, j) = \mathbf{x}_1 = \mathbf{R}_{11}^{-1} \mathbf{z}$ . This immediately implies that  $\tilde{\mathbf{R}}$  is still an upper triangular matrix, and only a  $j \times j$  triangular system needs be solved for  $\tilde{\mathbf{R}}(1 : j, j)$ . So the complexity of back substitution is reduced to  $\frac{1}{3}p^3$  FLOP.

An alternative yet less obvious algorithm is practically applied in numerical linear algebra software. Partition the triangular system as

$$\begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} \\ & \mathbf{R}_{22} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{R}}_{11} & \tilde{\mathbf{R}}_{12} \\ & \tilde{\mathbf{R}}_{22} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_{11} & \\ & \mathbf{I}_{22} \end{pmatrix},$$

then  $\tilde{\mathbf{R}}_{11} = \mathbf{R}_{11}^{-1}$ ,  $\tilde{\mathbf{R}}_{22} = \mathbf{R}_{22}^{-1}$  and  $\tilde{\mathbf{R}}_{12} = -\tilde{\mathbf{R}}_{11} \mathbf{R}_{12} \tilde{\mathbf{R}}_{22}$ . Let's call  $\mathbf{R}_{11}$  and  $\mathbf{R}_{22}$  as leading and trailing block, respectively. There is no restriction on the size of partition, so one special case is  $\mathbf{R}_{11} = \mathbf{R}(1 : (p-1), 1 : (p-1))$  and  $\mathbf{R}_{22} = \mathbf{R}(p, p)$ . However, for  $\mathbf{R}(1 : (p-1), 1 : (p-1))$ , the same kind of partition can be applied, and in this case the leading block has dimension  $(p-2) \times (p-2)$ . Such recursion can continue until the leading block becomes  $\mathbf{R}(1, 1)$ . This immediately implies an algorithm in Figure 4.3a. It has the same number of FLOP with the direct method, but it is rich in matrix-vector multiplication rather than back substitution, and is expected to be more efficient on modern computers. Croz and Higham (1992) covers both algorithms and their difference in numerical stability. Note that the algorithm can overwrite  $\mathbf{R}$  with  $\tilde{\mathbf{R}}$ . See Figure 4.2 for an illustration.



**Figure 4.3:** Algorithm for inverting a positive-definite matrix  $\mathbf{S}$  given its upper triangular Cholesky factor  $\mathbf{R}$ .

Transposed triangular matrix cross-product  $\mathbf{S}^{-1} = \tilde{\mathbf{R}}\tilde{\mathbf{R}}'$  is a special type of matrix multiplication. By exploiting that  $\mathbf{S}^{-1}$  is symmetric and  $\tilde{\mathbf{R}}$  is upper triangular, it is not hard to obtain an algorithm in Figure 4.3b for computing its upper triangular part. The algorithm also involves  $\frac{1}{3}p^3$  FLOP.

Now suppose we have a pivoted Cholesky factorization with pre-conditioning  $\mathbf{S} = \mathbf{R}^*\mathbf{R}^*$ , then  $\mathbf{R}^* = \mathbf{R}\mathbf{P}'\mathbf{J}$ , where  $\mathbf{R}$  is truly triangular. Computation of  $\mathbf{S}^{-1}$  in this case only needs a post-processing, since  $\mathbf{S}^{-1} = \mathbf{J}^{-1}\mathbf{P}\tilde{\mathbf{S}}\mathbf{P}'\mathbf{J}^{-1}$  and  $\tilde{\mathbf{S}} = \mathbf{R}^{-1}\mathbf{R}'^{-1}$  can be first formed using methods explained above.

#### 4.1.5 Symmetric eigen decomposition

Eigen decomposition is not just defined for a symmetric matrix, but since a symmetric  $p \times p$  matrix  $\mathbf{S}$  can be diagonalized, it gives a nice form of eigen decomposition:

$$\mathbf{S} = \mathbf{U}\mathbf{D}\mathbf{U}',$$

where  $\mathbf{U}$  is a  $p \times p$  orthonormal matrix containing eigenvectors, and  $\mathbf{D}$  is a diagonal matrix, whose diagonal elements are eigenvalues. This decomposition is often interpreted as an orthogonal reconstruction of  $\mathbf{S}$ , because

$$\mathbf{S} = \sum_{i=1}^p d_i \mathbf{u}_i \mathbf{u}_i',$$

where  $\mathbf{u}_i = \mathbf{U}(:, i)$  is the  $i^{\text{th}}$  eigen vector and  $d_i = \mathbf{D}(i, i)$  is the  $i^{\text{th}}$  eigen value. Since  $\|\mathbf{u}_i\| = 1$ ,  $d_i$  determines the magnitude of the contribution of  $d_i \mathbf{u}_i \mathbf{u}_i'$ . If we arrange eigenvalues so that  $|d_{i+1}| \geq |d_i|$ , then the contribution of these orthogonal components will be decreasing in magnitude. It is then natural to consider truncating the summation as an approximation to  $\mathbf{S}$ :

$$\mathbf{S}_k = \sum_{i=1}^k d_i \mathbf{u}_i \mathbf{u}_i' \approx \mathbf{S}, \quad k \leq p.$$

The larger  $k$  is, the better the approximation is. Such approximation is called a low-rank approximation, because if  $\mathbf{S}$  has a full rank of  $p$ ,  $\mathbf{S}_k$  only has a rank of  $k$ .

The above eigen decomposition implies that a matrix “root” of  $\mathbf{S}$  can be computed as  $\mathbf{S}^{\frac{1}{2}} = \mathbf{U}\mathbf{D}^{\frac{1}{2}}\mathbf{U}'$ , where  $\mathbf{D}^{\frac{1}{2}}$  is a diagonal matrix with  $\mathbf{D}^{\frac{1}{2}}(i, i) = \sqrt{d_i}$ , so that  $(\mathbf{S}^{\frac{1}{2}})^2 = \mathbf{S}^{\frac{1}{2}}\mathbf{S}^{\frac{1}{2}} = \mathbf{S}$ . When  $\mathbf{S}$  is positive-definite or positive-semidefinite, its root is still a real matrix. However, in this thesis,  $\mathbf{S}^{\frac{1}{2}}$  will be taken as the upper triangular Cholesky factor of  $\mathbf{S}$ , so that there is instead  $(\mathbf{S}^{\frac{1}{2}})'\mathbf{S}^{\frac{1}{2}} = \mathbf{S}$ .

Unlike any other matrix computation in this preliminary section, eigen decomposition is the only one with no exact FLOP count, because its practical computation is based on iterative methods. Precisely it is a two-stage task:

$$\begin{aligned}
v_1 &= \begin{bmatrix} \times \\ \times \\ \times \\ \times \end{bmatrix} \begin{bmatrix} \diamond & \times & \times & \times \\ \times & \bullet & \bullet & \bullet \\ \times & \bullet & \bullet & \bullet \\ \times & \bullet & \bullet & \bullet \end{bmatrix} \rightarrow \begin{bmatrix} \diamond & \times & \bullet & \bullet \\ \times & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{bmatrix} \\
v_2 &= \begin{bmatrix} \times \\ \times \\ \times \\ \times \end{bmatrix} \begin{bmatrix} \diamond & \diamond & \times & \times \\ \diamond & \times & \bullet & \bullet \\ \times & \bullet & \bullet & \bullet \\ \times & \bullet & \bullet & \bullet \end{bmatrix} \rightarrow \begin{bmatrix} \diamond & \diamond & \times & \bullet \\ \diamond & \times & \bullet & \bullet \\ \times & \bullet & \bullet & \bullet \\ \times & \bullet & \bullet & \bullet \end{bmatrix}
\end{aligned}$$

**Figure 4.4:** Illustration of Householder tridiagonalization using a  $4 \times 4$  matrix. See Figure 4.1 for meanings of place holders  $\bullet$ ,  $\times$  and  $\diamond$ .

**Table 4.1:** A brief summary of the three matrix factorizations elaborated in §4.1. “ $\sim$  upper triangular” means “equivalent to upper triangular”; that is,  $\mathbf{R}^*$  is upper triangular upon a column permutation.

	QR	Cholesky	symmetric eigen
math form	$\mathbf{X} = \mathbf{Q}\mathbf{R}^*$	$\mathbf{S} = \mathbf{R}^{*\prime}\mathbf{R}^*$	$\mathbf{S} = \mathbf{U}\mathbf{D}\mathbf{U}'$
input	$\mathbf{X}$ (general)	$\mathbf{S}$ (positive-semidefinite)	$\mathbf{S}$ (symmetric)
dimension	$n \times p$ ( $n \geq p$ )	$p \times p$	$p \times p$
output 1	$\mathbf{Q}$ (“factored form”)	$\mathbf{R}^*$ ( $\sim$ upper triangular)	$\mathbf{U}$ (orthogonal)
output 2	$\mathbf{R}^*$ ( $\sim$ upper triangular)		$\mathbf{D}$ (diagonal)
FLOP	$2np^2 - \frac{2}{3}p^3$	$\frac{1}{3}p^3$	$> \frac{4}{3}p^3$

1.  $\mathbf{S}$  is first tri-diagonalized by sequential Householder transformations from both left and right. This is similar to what happens in a QR factorization and Figure 4.4 is an illustration. Golub and Loan (2013, §8.3.1) gives a detailed algorithm for this, involving  $\frac{4}{3}p^3$  FLOP;
2. QR iteration is then applied on the tri-diagonal matrix to find eigenvalues. It will be too lengthy if I cover this topic here, but Golub and Loan (2013, §8.3) gives substantial details for anyone with interest. Each iteration has  $O(p^2)$  FLOP.

#### 4.1.6 Summary

This section is a consolidation of what has been covered so far. The three types of matrix factorizations are probably the most confusing to you, so I specially list them in Table 4.1. There is no implication that one factorization is more useful than the others; they are all essential in linear algebra. QR factorization is distinguished from the other two as it has almost no requirement on the input matrix. Symmetric eigen decomposition needs a symmetric matrix. Cholesky factorization has the strongest input requirement: all eigen values of the input symmetric matrix must be non-negative.

In practice, matrix factorizations can be used in various ways for different purposes. Here I will provide a few examples.

#### Solving a linear system of equations

They may be used for solving a linear equation  $\mathbf{A}\mathbf{x} = \mathbf{b}$ . QR factorization can be used to solve a very general system.

- For a full-rank system (exclusively when  $\mathbf{A}$  is square and has full rank) it is an alternative to Gaussian elimination or LU factorization;
- For an over-determined system (for example, when  $\mathbf{A}$  has more rows than columns) it yields the least squares solution;
- For an under-determined system (for example, when  $\mathbf{A}$  has more columns than rows) it constrains free parameters to zero and estimate a truncated full-rank system. This is also the way

that R function `lm` finds a least squares solution for rank-deficient system (i.e., when  $A$  has more rows than columns but it does not have full column rank).

Cholesky factorization can also be used for equation solving, but it only works for a positive-definite system, where  $A$  is square, symmetric and full-rank.  $A$  can be relaxed to be positive-semidefinite, in which case the system is essentially over-determined. Using Cholesky factorization with pivoting, numerical rank of the system can be estimated and a truncation can be performed, constraining free parameters to zeros hence yielding a least squares solution. Eigen decomposition can be used as an alternative to Cholesky factorization for solving this type of system, but in practice this is rarely done as it is more computationally costly.

Note that if the RHS of the system is zero, we get a homogenous system. The solution gives the *null space* of  $A$ . We will from time to time see a homogenous system. For example, find the eigenvectors for  $A$  associated with an eigenvalue  $\lambda$ . Mathematically eigenvectors are solutions to  $(A - \lambda I)x = 0$ . Of course, such is only used for educational purpose like in a textbook of undergraduate linear algebra. In practice, eigenvalues and eigenvectors are both unknown and can be found together using eigen decomposition (note: that I only introduced symmetric eigen decomposition does not mean eigen decomposition is not defined for general matrix; eigenvalues and eigenvectors are the most important character for a square matrix).

GAM estimation involves solving a penalized least squares problem (1.4). In fact, it is just a least squares problem and hence an equation solving problem. Consider a pivoted Cholesky factorization  $S_\lambda = E_\lambda^* E_\lambda$  (the Cholesky factor also depends on  $\lambda$ ), it is straightforward to rewrite the penalized least squares object as

$$\hat{\beta}_\lambda = \arg \min_{\beta} \left\| \begin{pmatrix} W^{\frac{1}{2}} y \\ 0 \end{pmatrix} - \begin{pmatrix} W^{\frac{1}{2}} X \\ E_\lambda \end{pmatrix} \beta \right\|^2, \quad (4.1)$$

which is a least squares problem with an augmented model matrix and an augmented response vector.

### Setting linear constraints on parameters

A linear system is also often used to specify linear constraints on parameters. This is mostly seen in constrained optimization problems. However, here in GAM computations I can give you a concrete example: centring constraint on a spline.

Since the individual  $f_i$  in an additive model is only estimable to within an intercept term, identifiability constraints need be applied. As discussed in Wood et al. (2013), a *sum-to-zero constraint* has the advantage of leading to narrow confidence intervals on the constrained  $f_i$ . For example, for  $f_1$  above this is  $\sum_{i=1}^n f_1(x_i) = 0$  and for  $f_2$  this is  $\sum_{i=1}^n f_2(z_{i1}, z_{i2}) = 0$ . Such constraint can also be written as a matrix equation  $C_i' X_i = 0$ , where  $C_i = X_i' \mathbf{1}$  is a vector containing column-wise sum of  $X_i$ .

It is easy to reparametrize  $\beta_i$  to incorporate the constraint directly into  $X_i$  and  $S_i$ . The constraint implies that  $\beta_i$  lies in the null space of  $C_i'$ . Then we can solve the homogenous system  $C_i' X_i = 0$  for the null space basis  $B_i$ . Clearly there exists some unconstrained parameters  $\tilde{\beta}_i$  such that  $\beta_i = B_i \tilde{\beta}_i$ . So instead of estimating  $\beta_i$  we can estimate  $\tilde{\beta}_i$ . This is essentially a reparametrization so that we have a new design matrix  $X_i B_i$  and new coefficient vector  $\tilde{\beta}_i$  for spline  $f_i$ .

Note that the QR factorization provides another way for finding  $B_i$  than equation solving. As mentioned in §4.1.1, following a thin QR factorization  $A = QR$ ,  $Q^\perp$  is an orthonormal basis for the null space of  $A'$ . Therefore, to find an orthonormal basis for  $A$  we just need to apply QR factorization to  $A'$ . So, following a thin QR factorization  $C_i = Q_i R_i$ , spline  $f_i$  can be reparametrized using  $B_i = Q_i^\perp$ , ending up with a new design matrix  $X_i Q_i^\perp$ . As has been explained in §4.1.1, matrix multiplication with  $Q_i^\perp$  is effectively a set of Householder transformations. In particular, as  $C_i$  is a

single-column matrix, there will be just one Householder transformation here. The new design matrix is one column less, and the new penalty matrix will lose a row and a column. This reparametrization is said to “absorb centring constraint into  $\mathbf{f}_i$ ”.

## Variable transformation

Absorbing centring constraint is an example of variable transformation using QR factorization. Now I will give other examples for Cholesky factorization and eigen decomposition.

In §2.2.2 I have defined the standardization for an AR(1) process. Broadly speaking this provides a way to make transformations between variables with a general covariance matrix and variables with an identity covariance matrix. That is, we can sample from a Gaussian distribution with a general covariance matrix by transforming i.i.d. Gaussian samples, and vice versa.

- Suppose we want to sample  $\mathbf{x} \sim N(\mathbf{0}, \mathbf{V}\phi)$ . Consider a Cholesky factorization  $\mathbf{V} = \mathbf{L}\mathbf{L}'$  where  $\mathbf{L}$  is the lower triangular Cholesky factor.  $\mathbf{x}$  can be obtained by first sampling  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}\phi)$  then setting  $\mathbf{x} = \mathbf{L}\mathbf{z}$ . (The reverse statement is:  $\mathbf{x}$  can be standardized via  $\mathbf{L}^{-1}\mathbf{x}$ .)
- Suppose we want to sample  $\mathbf{x} \sim N(\mathbf{0}, \mathbf{W}^{-1}\phi)$ . Consider a Cholesky factorization  $\mathbf{W} = \mathbf{U}'\mathbf{U}$  where  $\mathbf{U}$  is the upper triangular Cholesky factor.  $\mathbf{x}$  can be obtained by first sampling  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}\phi)$  then setting  $\mathbf{x} = \mathbf{U}^{-1}\mathbf{z}$ . (The reverse statement is:  $\mathbf{x}$  can be standardized via  $\mathbf{U}^{-1}\mathbf{x}$ .)

Eigen decomposition can be used for the same task. In the first case, decompose the correlation matrix  $\mathbf{V} = \mathbf{E}\mathbf{D}\mathbf{E}'$ , then we have  $\mathbf{L} = \mathbf{E}\mathbf{D}^{\frac{1}{2}}$ . In the second case, decompose the inverse correlation matrix  $\mathbf{W} = \mathbf{E}\mathbf{D}\mathbf{E}'$ , then we have  $\mathbf{U} = \mathbf{D}^{\frac{1}{2}}\mathbf{E}'$ . (Note that  $\mathbf{L}$  and  $\mathbf{U}$  are not triangular matrices in this method.)

## Computing determinant of a matrix

Matrix factorizations provide a way to compute determinant of a square matrix.

- For a general square matrix, LU factorization is used:  $\mathbf{A} = \mathbf{L}\mathbf{U}$ , where  $\mathbf{L}$  is a unit lower triangular matrix (a lower triangular matrix with all ones on its main diagonal) and  $\mathbf{U}$  is an upper triangular matrix. Then,  $|\mathbf{A}| = |\mathbf{L}||\mathbf{U}| = |\mathbf{U}|$ . (I haven’t introduced LU factorization in the previous sections, so there is no need to know anything further than this).
- For a general symmetric matrix, symmetric eigen factorization can be used:  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}'$ , then  $|\mathbf{A}| = |\mathbf{U}||\mathbf{D}||\mathbf{U}'| = |\mathbf{U}|^2|\mathbf{D}| = |\mathbf{D}| = \prod_i \mathbf{D}(i, i)$ .
- For a positive-definite matrix, Cholesky factorization can be used:  $\mathbf{A} = \mathbf{R}'\mathbf{R}$ , then  $|\mathbf{A}| = |\mathbf{R}||\mathbf{R}'| = |\mathbf{R}|^2 = \prod_i \mathbf{R}(i, i)^2$ .

GAM estimation also involves computation of determinant (or actually log-determinant). If you still remember (check (1.9)), there are three determinant terms in the REML score  $\mathcal{V}_r$  (revisit §1.1.6 if you need to):  $\log|\mathbf{S}_\lambda|_+$ ,  $\log|\mathbf{H}_\lambda|$  and  $\log|\mathbf{W}|$ . The last determinant is mostly trivial. Very often  $\mathbf{W}$  is just a diagonal matrix of weights  $w_i$ , so the determinant is just  $|\mathbf{W}| = \prod_i w_i$  or  $\log|\mathbf{W}| = \sum_i \log(w_i)$ . In a slightly more complicated case  $\mathbf{W}$  is symmetric tri-diagonal coming from AR(1) model errors, and in §2.2.2 I have already shown what this determinant is. In the forthcoming sections of this Chapter, you will see how  $\log|\mathbf{S}_\lambda|_+$  and  $\log|\mathbf{H}_\lambda|$  are computed using matrix factorizations.

## Achieving computational efficiency

Generally if you do some similar matrix computations in a loop where one matrix is loop invariant, there might be a change to speed up the loop by an appropriate initial factorization of this matrix.

That matrix factorization helps achieve computational efficiency will be seen a few occasions in the rest of this Chapter, like the model matrix reduction in §4.2 and the initial reparametrization for REML estimation in §4.3.1. Here I will provide another example.

Let  $\mathbf{A}$  be a  $p \times p$  real symmetric matrix. How would you compute a matrix power series  $\sum_{i=1}^m \mathbf{A}^i$ ? The obvious yet the worst idea is to do the summation via a loop, computing  $\mathbf{A}^i$  by  $(i-1)$  matrix multiplication  $\prod_{j=1}^i \mathbf{A}_j$  in the  $i^{\text{th}}$  iteration. This involves  $O(m^2 p^3)$  FLOP. A better arrangement is to accumulate  $\mathbf{B}_i = \prod_{j=1}^i \mathbf{A}_j$  along the loop so that only one matrix multiplication  $\mathbf{B}_{i+1} = \mathbf{B}_i \mathbf{A}$  is needed in each iteration. This reduces the computational complexity to  $O(mp^3)$ .

Now consider a symmetric eigen decomposition  $\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{U}'$ , then we see  $\sum_{i=1}^m \mathbf{A}^i = \sum_{i=1}^m (\mathbf{U} \mathbf{D} \mathbf{U}') = \mathbf{U} (\sum_{i=1}^m \mathbf{D}^i) \mathbf{U}'$ . Note that  $\mathbf{D}$  is a diagonal matrix with diagonal elements  $d_1, d_2, \dots, d_p$ , so that its power series is just

$$\sum_{i=1}^m \begin{pmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_p \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^m d_1 & & & \\ & \sum_{i=1}^m d_2 & & \\ & & \ddots & \\ & & & \sum_{i=1}^m d_p \end{pmatrix},$$

where the  $j^{\text{th}}$  element can be analytically computed by  $\sum_{i=1}^m d_j = \frac{d_j(1-d_j^m)}{1-d_j}$  without a loop. So the computational complexity of this method is really just the  $O(p^3)$  FLOP for the initial decomposition; the summation loop has no costs at all!

## Conclusion

I have presented rich contents here on matrix computations. While they are lengthy, they are right on the theme of this Chapter (and this thesis!). You will shortly see many matrix operations and if these operations are new to you, you can easily be buried in details. In this situation, do remind yourself of the aim the matrix computations summarized in the headings for the above paragraphs.

## 4.2 Model matrix reduction for additive models

The core of estimating an additive model is just solving a penalized least squares (1.4) or (4.1). And if the smoothing parameters  $\boldsymbol{\lambda}$  are unknown, the problem is iteratively solved many times until the best  $\boldsymbol{\lambda}$  is found that minimizes, say, the REML score  $\mathcal{V}_r$ . The most obvious way to do this iterative search is a grid search. However, I have already stated in §2.2.2 that a grid search is much too inefficient. If the parameter space has high dimension, this simple method would be prohibitively costly (as the number of grid points required grows geometrically with dimension; this is known as the “curse of dimensionality”). Generally a guided search is performed. For example, the golden-section search in §2.2.3 is just a simple guided search for 1D optimization that converges fast (enough). Minimization of the REML score  $\mathcal{V}_r$  is a high dimensional optimization (as there can be many splines and smoothing parameters), and Newton-Raphson algorithm is a guided search that converges fast (enough).

More details on Newton-Raphson algorithm will be provided in the next section (§4.3). For the moment, let us focus on that fact that (4.1) needs be solved many times in a loop. Remember the “Achieving computational efficiency” headings in the previous section? Obviously  $\mathbf{X}$  is loop-invariant and there might be an opportunity to apply matrix factorization here. For additive models with known  $\mathbf{W}$ ,  $\mathbf{W}^{\frac{1}{2}} \mathbf{X}$  will be loop-invariant. Consider a pivoted QR factorization (see §4.1.1 if you need a revision)  $\mathbf{W}^{\frac{1}{2}} \mathbf{X} = \mathbf{Q} \mathbf{R}^*$ , then  $D_p(\boldsymbol{\beta})$  in (1.5) can be rewritten as

$$D_p(\boldsymbol{\beta}, \boldsymbol{\lambda}) = \|\mathbf{f} - \mathbf{R}^* \boldsymbol{\beta}\|^2 + \boldsymbol{\beta}' \mathbf{S}_{\boldsymbol{\lambda}} \boldsymbol{\beta} + \|\mathbf{W}^{\frac{1}{2}} \mathbf{y}\|^2 - \|\mathbf{f}\|^2, \quad (4.2)$$

where  $\mathbf{f} = \mathbf{Q}'\mathbf{W}^{\frac{1}{2}}\mathbf{X}$ . Note that  $\|\mathbf{W}^{\frac{1}{2}}\mathbf{y}\|^2 - \|\mathbf{f}\|^2$  is also loop-invariant; so what really needs be solved in a loop is a reduced penalized least squares problem

$$\|\mathbf{f} - \mathbf{R}^*\boldsymbol{\beta}\|^2 + \boldsymbol{\beta}'\mathbf{S}_\lambda\boldsymbol{\beta}.$$

Although I haven't told you what the computational costs associated with a penalized least squares problem is, you should have observed that matrices in the reduced problem have smaller size than those in the original version (1.4) or (4.1), hence been able to infer that solving the reduced version in a loop is much cheaper.

In fact, for large datasets (i.e., large  $n$ ) this is double gains. The model matrix  $\mathbf{X}$  has large number of rows which can take huge amount of RAM for storage. For example, a  $300000 \times 10000$  matrix takes 22.35 GB. The reduced penalized least squares objective eliminates the need to access it every iteration hence the loop will be more memory efficient.

Such initial QR factorization for the model matrix  $\mathbf{X}$  to obtain  $\mathbf{f}$  and  $\mathbf{R}^*$  prior to  $\mathcal{V}_r$  minimization, is termed *model matrix reduction*, or precisely the *QR reduction* variant in this thesis. There is an alternative variant through matrix cross-product  $\mathbf{X}'\mathbf{W}\mathbf{X}$  and a subsequent Cholesky factorization, termed *pseudo QR reduction*. This section will give computational details for both variants.

You probably wonder why this initial QR factorization worths writing a section, given that QR factorization has been well explained in §4.1.1. Well, the issue is that  $\mathbf{X}$  is so large. Even if it won't be accessed in the REML estimation loop, a standard QR factorization still requires this large matrix to be in memory in the first place. Holding this matrix in RAM is challenging or even impossible on most computational devices. Particularly, there is no upper bound for  $n$ , which means there is no upper bound for memory usage in this standard way. The sections on model matrix reduction here aim to demonstrate an implementation where memory requirement is a constant that is independent of  $n$ . It will also turn out that such implementation can exploit parallel computing for speed.

To simplify notation in the following demonstration, I will assume that the model matrix and the response vector have absorbed weights:  $\mathbf{X} \leftarrow \mathbf{W}^{\frac{1}{2}}\mathbf{X}$ ,  $\mathbf{y} \leftarrow \mathbf{W}^{\frac{1}{2}}\mathbf{y}$ .

#### 4.2.1 QR reduction

In fact, the solution to avoid forming  $\mathbf{X}$  as a whole has been established, by using a QR updating algorithm for iterative update of a QR factorization(Sven and Craig, 2008, §2.3.2).

The QR updating problem is an example of many “online” problems. Here a simple example to help understand what an “online” problem is. Suppose we want to computed and report sample variance of a dataset where a new datum arrives every second. Assume that the dataset initially has a single datum at second 1, then it will contain  $t$  data by second  $t$  and we want to report their sample variance  $v_t$ . A poor man's approach would be using formula  $m_t = \frac{1}{t} \sum_{i=1}^t x_i$ ,  $v_t = \sum_{i=1}^t (x_i - m_i)^2$ , which scans every record up to time  $t$  thus involves  $O(t)$  computation complexity. This implies that computation is increasingly time-consuming as time goes by, and eventually we won't be able to output  $v_t$  as fast as  $x_t$  arrives. Do we have to scan the data vector right from the beginning every second? Is there a way to simply update  $v_{t-1}$  to get  $v_t$ ? Yes, there is. Starting with  $m_1 = x_1$  and  $s_1 = 0$ , the Welford's online algorithm(Welford, 1962) suggests  $m_t = m_{t-1} + (x_t - m_{t-1}) / t$ ,  $s_t = s_{t-1} + (x_t - m_{t-1})(x_t - m_t)$ ,  $v_t = s_t / t$  which updates the sample mean  $m_t$  and sample variance  $v_t$  at  $O(1)$  computational complexity for any  $t \geq 2$ . Such algorithm is called “online” because it does not visit or store old records: we may discard  $x_{t-1}$  as soon as  $m_t$  and  $v_t$  are obtained.

Now let us get back to the QR factorization problem. Partition the large model matrix  $\mathbf{X}$  and



response vector  $\mathbf{y}$  into row chunks:

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_i \\ \vdots \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_i \\ \vdots \end{pmatrix},$$

and imagine those chunks arrive by time. There are no restrictions on the size of each chunk, but to ease practical implementation let all chunks have equal row chunk size  $c$  (although the final chunk may have fewer rows if  $c$  does not divide  $n$ ). Denote  $\mathbf{X}^{[i]}$  and  $\mathbf{y}^{[i]}$  as the previous  $i$  row chunks of  $\mathbf{X}$  and  $\mathbf{y}$ , then there are the following update formulas.

1. Given a pivoted QR factorization  $\mathbf{X}^{[i]} = \mathbf{Q}_i \mathbf{R}_i^*$ , the pivoted QR factorization for  $\mathbf{X}^{[i+1]}$  is only up to a factorization of  $\begin{pmatrix} \mathbf{R}_i^* \\ \mathbf{X}_{i+1} \end{pmatrix}$ , since

$$\mathbf{X}^{[i+1]} = \begin{pmatrix} \mathbf{X}^{[i]} \\ \mathbf{X}_{i+1} \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_i \mathbf{R}_i^* \\ \mathbf{X}_{i+1} \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_i & \\ & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{R}_i^* \\ \mathbf{X}_{i+1} \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_i & \\ & \mathbf{I} \end{pmatrix} \tilde{\mathbf{Q}}_{i+1} \mathbf{R}_{i+1}^*.$$

So  $\mathbf{Q}_{i+1} = \begin{pmatrix} \mathbf{Q}_i & \\ & \mathbf{I} \end{pmatrix} \tilde{\mathbf{Q}}_{i+1}$  and  $\mathbf{R}_{i+1}^*$  are ‘Q’ and ‘R’ factors of  $\mathbf{X}^{[i+1]}$ .

2. Given  $\mathbf{f}_i = \mathbf{Q}_i' \mathbf{y}^{[i]}$ ,  $\mathbf{f}_{i+1} = \mathbf{Q}_{i+1}' \mathbf{y}^{[i+1]}$  can be obtained by an update

$$\mathbf{f} = \mathbf{Q}_{i+1}' \mathbf{y}^{[i+1]} = \tilde{\mathbf{Q}}_{i+1}' \begin{pmatrix} \mathbf{Q}_i' & \\ & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{y}^{[i]} \\ \mathbf{y}_i \end{pmatrix} = \tilde{\mathbf{Q}}_{i+1}' \begin{pmatrix} \mathbf{f}_i \\ \mathbf{y}_i \end{pmatrix}.$$

In other words, there is no realistic need to form  $\mathbf{Q}_{i+1}$ .

The above is an “online” algorithm, because

1.  $\mathbf{R}_{i+1}^*$  and  $\mathbf{f}_{i+1}$  are updated from  $\mathbf{X}_i$ ,  $\mathbf{R}_i^*$  and  $\mathbf{f}_i$  rather than computed from  $\mathbf{X}^{[i+1]}$  and  $\mathbf{y}^{[i+1]}$ ,
2. as soon as  $\mathbf{R}_{i+1}$  and  $\mathbf{f}_{i+1}$  are obtained  $\mathbf{X}_i$  can be discarded.

As a result, only a single row chunk of  $\mathbf{X}$  and  $\mathbf{y}$  need be in RAM for this “online updating”, eliminating the storage requirement for the whole  $\mathbf{X}$ . For later reference, let us call the basic implementation where the whole  $\mathbf{X}$  is formed the “offline” implementation.

Figure 4.5 gives a practical implementation of the “online” algorithm. In each iteration, only a  $\tilde{c} \times p$  row chunk of  $\mathbf{X}$  is formed, and only  $O(cp + p^2)$  storage is needed in RAM for the QR factorization.

It is also easy to count that the algorithm involves  $(2np^2 - \frac{2}{3}p^3 + 2(\frac{n}{c} - 1)p^3)$  FLOP. Compared with the  $(2np^2 - \frac{2}{3}p^3)$  FLOP in the “offline” version, the algorithm has a  $2(\frac{n}{c} - 1)p^3$  overhead, which is approximately  $\frac{p}{c}$  of the  $(2np^2 - \frac{2}{3}p^3)$  useful work when  $n \gg p$ . To keep this overhead low, it makes sense to choose  $c$  as large as possible. But under a memory constraint  $cp < M$ , we can only set  $c = \frac{M}{p}$  at best, thus the percentage of overhead becomes  $p^2/M$ . This implies that on a machine with a certain amount of RAM, the overhead is negligible for small  $p$ , but has a quadratic growth for increasingly larger  $p$ . This is a major drawback of the algorithm.

#### 4.2.2 Pseudo QR reduction

I will now introduce the other variant for model matrix reduction: pseudo QR reduction. It is motivated by the link between Cholesky factorization and QR factorization, as described in §4.1.2.

---

initialize  $\mathbf{R}^*$  as a  $0 \times p$  matrix and  $\mathbf{f}$  as a length-0 vector

$\mathbf{R}^* = \text{zeros}(0, p)$

$\mathbf{f} = \text{zeros}(0)$

$i = 1$

while ( $i \leq n$ )

    determine the row chunk size for the current iteration

$\tilde{c} = \min\{c, n - i + 1\}$

    form a row chunk of  $\mathbf{X}$  and  $\mathbf{y}$

$\tilde{\mathbf{X}} = \mathbf{X}(i : (i + \tilde{c} - 1), :)$

$\tilde{\mathbf{y}} = \mathbf{y}(i : (i + \tilde{c} - 1))$

    QR factorization and orthonormal transformation

$\begin{pmatrix} \mathbf{R}^* \\ \tilde{\mathbf{X}} \end{pmatrix} = \tilde{\mathbf{Q}} \tilde{\mathbf{R}}$

$\mathbf{f} = \tilde{\mathbf{Q}}' \begin{pmatrix} \mathbf{f} \\ \tilde{\mathbf{y}} \end{pmatrix}$

    update  $\mathbf{R}^*$  and  $\mathbf{f}$

$\mathbf{R}^* \leftarrow \tilde{\mathbf{R}}$

$\mathbf{f} \leftarrow \tilde{\mathbf{f}}$

    move to next row chunk

$i += \tilde{c}$

---

**Figure 4.5:** “Online” QR reduction with a row chunk size  $c$ . By forming a row chunk of  $\mathbf{X}$  per iteration, the memory footprint is bounded by  $O(cp + p^2)$ . The QR factorization in each iteration involves  $(2\tilde{c}p^2 + \frac{4}{3}p^3)$  FLOP, and the algorithm involves  $(2np^2 - \frac{2}{3}p^3 + 2(\frac{n}{c} - 1)p^3)$  FLOP.

By first forming a matrix cross-product  $\mathbf{X}'\mathbf{X}$ , a subsequent pivoted Cholesky factorization (with pre-conditioning for computational stability) would yield  $\mathbf{R}^*$ . However,  $\mathbf{f}$  is not available without a ‘Q’ factor, but the QR factorization implies  $\mathbf{Q} = \mathbf{X}\mathbf{R}^{*-1}$ , hence we have  $\mathbf{f} = \mathbf{Q}'\mathbf{y} = \mathbf{R}^{*-1}(\mathbf{X}'\mathbf{y})$ , which can be computed by solving a triangular system of linear equations (see §4.1.3 if you need a revision). Note that when  $\mathbf{X}$  is rank-deficient,  $\mathbf{R}^*$  is not invertible. Computation of  $\mathbf{f}$  then proceeds as follows. Recall from §4.1.2 that we have  $\mathbf{R}^* = \mathbf{R}\mathbf{P}'\mathbf{J}$  where  $\mathbf{R}$  is an upper triangular matrix. Let  $r$  be the rank of  $\mathbf{R}$ , then we first initialize  $\mathbf{f}$  as a length- $p$  vector of zeros, then set  $\mathbf{f}(1:r) = \mathbf{R}(1:r, 1:r)^{-1}((\mathbf{P}'\mathbf{J}^{-1}\mathbf{X}'\mathbf{y})(1:r))$ .

A big advantage of pseudo QR reduction over QR reduction is its efficiency. It involves  $np^2$  FLOP for the cross-product,  $\frac{1}{3}p^3$  FLOP for the Cholesky factorization and  $p^2$  FLOP for solving the triangular system, hence a total of  $(np^2 + \frac{1}{3}p^3 + p^2)$  FLOP, approximately half of the  $(2np^2 - \frac{2}{3}p^3)$  FLOP with QR reduction when  $n \gg p$ .

Pseudo QR reduction also has an “online” implementation. After row chunking of  $\mathbf{X}$  and  $\mathbf{y}$ , there are  $\mathbf{X}'\mathbf{X} = \sum_i \mathbf{X}'_i \mathbf{X}_i$  and  $\mathbf{X}'\mathbf{y} = \sum_i \mathbf{X}'_i \mathbf{y}_i$ , that is, we form one row chunk in each iteration and accumulate  $\mathbf{X}'\mathbf{X}$  and  $\mathbf{X}'\mathbf{y}$ . Figure 4.6 is a sketch of its practical implementation. A key message is that the algorithm has no chunking overhead compared with the non-iterative version.

### 4.2.3 Parallel computing

The “online” algorithms for both reduction methods remove memory bottleneck, but do not improve the fundamental  $O(np^2)$  computational complexity. So it is expected that as  $n$  grows, computation becomes more time consuming.

A practical idea to speed up the computation is *parallel computing*. For example, if a loop runs 100 iterations, we may do the first 50 iterations on one processor and the other 50 iterations on a second processor. Provided that the latter 50 iterations do not rely on the processing results of the former 50 iterations, we can execute the two parts at the same time (or in parallel), so that while the total amount of work remains unchanged, the realistic execution time may be reduced (ideally) by half.

---

```

initialize a  $p \times p$  matrix of zeros and a length- $p$  vector of zeros
 $\mathbf{A} = \text{zeros}(p, p)$ 
 $\mathbf{z} = \text{zeros}(p)$ 
 $i = 1$ 
while ( $i \leq n$ )
    determine the row chunk size for the current iteration
     $\tilde{c} = \min\{c, n - i + 1\}$ 
    form a row chunk of  $\mathbf{X}$  and  $\mathbf{y}$ 
     $\tilde{\mathbf{X}} = \mathbf{X}(i : (i + \tilde{c} - 1), :)$ 
     $\tilde{\mathbf{y}} = \mathbf{y}(i : (i + \tilde{c} - 1))$ 
    accumulate matrix-matrix cross-product and matrix-vector cross-product
     $\mathbf{A} += \tilde{\mathbf{X}}' \tilde{\mathbf{X}}$ 
     $\mathbf{z} += \tilde{\mathbf{X}}' \tilde{\mathbf{y}}$ 
    move to next row chunk
     $i += \tilde{c}$ 
Cholesky factorization and solving a triangular system
 $\mathbf{A} = \mathbf{R}^{*'} \mathbf{R}^*$ 
 $\mathbf{f} = \mathbf{R}^{*'}{}^{-1} \mathbf{z}$ 

```

---

**Figure 4.6:** “Online” pseudo QR reduction with a row chunk size  $c$ . By forming a row chunk of  $\mathbf{X}$  per iteration, the memory footprint is bounded by  $O(cp)$ . The cross-product in each iteration involves  $\tilde{c}p^2$  FLOP, the final Cholesky factorization involves  $\frac{1}{3}p^3$  FLOP and solving the triangular system involves  $p^2$  FLOP. The algorithm involves a total of  $(np^2 + \frac{1}{3}p^3 + p^2)$  FLOP.

For later reference, computation performed strictly on a single processor is called *serial computing*.

It is probably not obvious how algorithms in Figures 4.5 and 4.6 can be made parallel, as there is loop carried dependence: the  $(i+1)^{\text{th}}$  iteration relies on the execution result of the  $i^{\text{th}}$  iteration. However, both algorithms can be adapted to break such dependency. Suppose there are  $m > 1$  processors. We split  $\mathbf{X}$  and  $\mathbf{y}$  into  $m$  chunks:

$$\begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_m \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_m \end{pmatrix}.$$

For QR reduction method, let the  $j^{\text{th}}$  processor perform  $\mathbf{X}_j = \mathbf{Q}_j \mathbf{R}_j^*$  and  $\mathbf{f}_j = \mathbf{Q}_j' \mathbf{y}_j$  using the “online” algorithm in Figure 4.5, then the following identity implies that processing results  $\mathbf{R}_j$ ’s from all processors can be stacked for a final QR factorization to produce the ‘Q’ and ‘R’ factors associated with the complete  $\mathbf{X}$ .

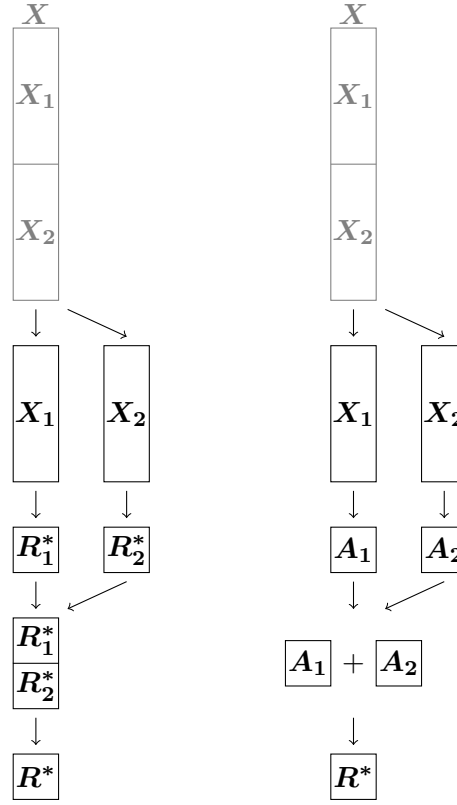
$$\begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_m \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_1 & & & \\ & \mathbf{Q}_2 & & \\ & & \ddots & \\ & & & \mathbf{Q}_m \end{pmatrix} \begin{pmatrix} \mathbf{R}_1^* \\ \mathbf{R}_2^* \\ \vdots \\ \mathbf{R}_m^* \end{pmatrix} = \begin{pmatrix} \mathbf{Q}_1 & & & \\ & \mathbf{Q}_2 & & \\ & & \ddots & \\ & & & \mathbf{Q}_m \end{pmatrix} \mathbf{Q} \mathbf{R}^*.$$

It also follows that

$$\mathbf{f} = \mathbf{Q}' \begin{pmatrix} \mathbf{Q}_1 & & & \\ & \mathbf{Q}_2 & & \\ & & \ddots & \\ & & & \mathbf{Q}_m \end{pmatrix} \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_m \end{pmatrix} = \mathbf{Q}' \begin{pmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \\ \vdots \\ \mathbf{f}_m \end{pmatrix}.$$

For pseudo QR reduction method, let the  $j^{\text{th}}$  processor perform  $\mathbf{A}_j = \mathbf{X}_j' \mathbf{X}_j$  and  $\mathbf{z}_j = \mathbf{X}_j' \mathbf{y}_j$  using the “online” algorithm in Figure 4.6, then the processing results from all processors can be accumulated by  $\mathbf{A} = \sum_{j=1}^m \mathbf{A}_j$ ,  $\mathbf{z} = \sum_{j=1}^m \mathbf{z}_j$  to get the results associated with the complete  $\mathbf{X}$  for the Cholesky factorization and equation solving.

Such parallelism is called “fork-join” (see Figure 4.7 for an illustration with  $m = 2$ ; note that  $\mathbf{y}$



**Figure 4.7:** Illustration of the parallel model matrix reduction implemented in `bam`. Assuming that there are two processors, the initial model matrix  $\mathbf{X}$  is split evenly into two row chunks  $\mathbf{X}_1$  and  $\mathbf{X}_2$  for processor 1 and processor 2. In QR reduction, illustrated on the left flow, a QR factorization  $\mathbf{X}_j = \mathbf{Q}_j \mathbf{R}_j^*$  is performed in parallel on two processors, then  $\mathbf{R}_1^*$  and  $\mathbf{R}_2^*$  are stacked by row for a last QR factorization for the final  $\mathbf{R}^*$ . In pseudo QR reduction, illustrated in the right flow, a matrix cross-product  $\mathbf{A}_j = \mathbf{X}_j' \mathbf{X}_j$  is performed on each processor, then results are accumulated  $\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2$  for a final Cholesky factorization  $\mathbf{A} = \mathbf{R}^{*'} \mathbf{R}^*$  to obtain  $\mathbf{R}^*$ . Computations on each processor can further use “online” algorithms in Figures 4.5 and 4.6 to keep memory footprint under control. Note that  $\mathbf{X}$  is drawn and annotated in gray, because it is not explicitly constructed. The initial splitting and dispatching only work on row index of  $\mathbf{X}$ .

is not sketched in the Figure). It is ideal on a high-performance computing cluster, where a large model matrix can be split and dispatched to different computing nodes. The `bam` function supports such parallel processing, but since R software is built on a shared-memory not distributed-memory architecture, it treats each CPU core of a multi-core computer as a processor.

The chunking applied to  $\mathbf{X}$  and  $\mathbf{y}$  here is not to be confused with the chunking done in “online” algorithms; in practice the latter is nested in the former. Also, rows of  $\mathbf{X}$  and  $\mathbf{y}$  should be split roughly evenly between processors for *load balancing*. With near equal workload on each processor, all processors can finish their share at about the same time; otherwise the execution time depends on the processor that finishes its work last.

The “fork” step is where parallel computing is truly profitable; the “join” step is actually an overhead since the work is done by a single processor. Suppose that  $m > 1$  processors are used at the “fork” step, we have the following observation.

- The “join” step of QR reduction factorizes an  $mp \times p$  matrix. This involves  $(2mp^3 - \frac{2}{3}p^3)$  FLOP which grows linearly with  $m$ .
- The “join” step of pseudo QR reduction adds up  $(m-1) p \times p$  matrices and performs a Cholesky factorization and equation solving. This involves  $((m-1)p^2 + \frac{1}{3}p^3 + p^2)$  FLOP, which is almost constant in  $m$  for even moderately big  $p$  (since the leading order term  $\frac{1}{3}p^3$  does not depend on  $m$ ).

Thus one may expect that as  $m$  increases, the practical performance of parallel QR reduction would be bottlenecked by the increasingly expensive “join” step.

We will get more insight if we examine *parallel speedup* or *parallel scaling factor*. Practical parallel speedup is defined with realistic execution time by

$$\text{practical parallel scaling factor} = \frac{\text{computation time using 1 core}}{\text{computation time using } m \text{ cores}}.$$

It is a measure of the effectiveness of parallel computing. Theoretical speedup is calculated with computational complexity by

$$\text{theoretical parallel scaling factor} = \frac{\text{FLOP count in serial computing}}{\text{FLOP count per core} + \text{FLOP count at “join”}}.$$

It is an upper bound for practical parallel speedup. For QR reduction this quantity is

$$\frac{2np^2 - \frac{2}{3}p^3 + 2(\frac{n}{c} - 1)p^3}{(2\frac{n}{m}p^2 - \frac{2}{3}p^3 + 2(\frac{n}{mc} - 1)p^3) + (2mp^3 - \frac{2}{3}p^3)} \approx \frac{m}{1 + \frac{m^2pc}{n(p+c)}},$$

and for pseudo QR reduction it is

$$\frac{np^2 + \frac{1}{3}p^3 + p^2}{\frac{n}{m}p^2 + ((m-1)p^2 + \frac{1}{3}p^3 + p^2)} \approx \frac{m}{1 + \frac{(m-1)p}{3n+p}}.$$

These mathematical results can be summarized as follows.

- Both values are below  $m$  so perfect scaling is never possible, i.e., using  $m$  cores never reduces the execution time to  $1/m$  of that in serial computing.
- The  $m^2$  in the result for QR reduction implies that the parallel scaling for this method is likely to degrade fast as  $m$  increases. Note that having  $n$  in the denominator does not mean that the fraction would be small. The value also depends on  $p$  and  $c$ . For example, consider  $n = 10^6$ ,  $p = 5 \times 10^3$  and  $c = 2 \times 10^4$ , then  $\frac{pc}{n(p+c)} \approx 0.0046$ . When multiplied by  $m^2$ , say  $8^2 = 64$ , it gives nearly 0.3.
- The  $\frac{(m-1)p}{3n+p}$  in the result for pseudo QR reduction also hints a performance degrade as  $m$  increases. However, it tends to have a negligible impact in practice if  $n \gg p$ , or even if  $n$  is just moderately bigger than  $p$ . So it is predicted that this method should have very good practical parallel scaling.

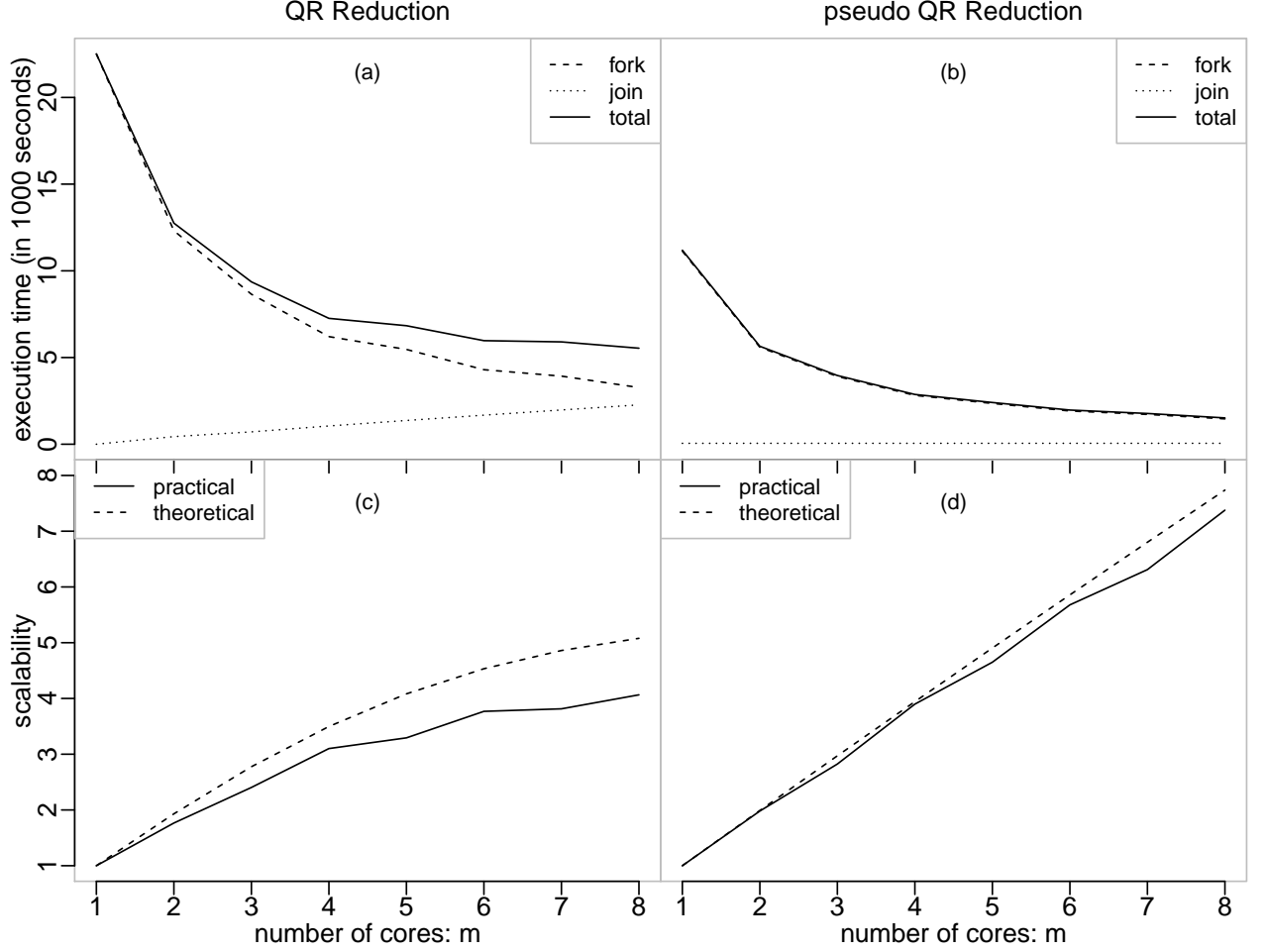
See Figure 4.8 for a practical experiment verifying those findings.

#### 4.2.4 Summary

Below is a comparison between QR reduction and pseudo QR reduction consolidating previous demonstrations.

- QR reduction and pseudo QR reduction respectively approximately involve  $2np^2$  and  $np^2$  FLOP, thus the latter is practically 100% times faster (verified by the timing statistics for two methods in Figure 4.8 when  $m = 1$ );
- In “online” implementation, QR reduction has a chunking overhead that grows at a rate of  $O(p^2)$ , while pseudo QR reduction has no chunking overhead at all;

computation time (in seconds) for “fork-join” parallel model matrix reduction									
number of cores: $m$		1	2	3	4	5	6	7	8
QR reduction	fork	22507	12310	8652	6205	5466	4298	3927	3269
	join	0	439	712	1053	1369	1675	1975	2268
	total	22507	12749	9363	7258	6835	5972	5902	5537
pseudo QR reduction	fork	11136	5596	3910	2817	2353	1917	1720	1463
	join	52	53	53	53	53	53	53	54
	total	11188	5649	3963	2870	2406	1970	1773	1517



**Figure 4.8:** Practical performance of “fork-join” parallel QR reduction and pseudo QR reduction. Experiment is undertaken on an Intel Xeon E5-2650 v2 workstation exploiting up to 8 CPU cores (see Appendix A for hardware information). An arbitrary  $430000 \times 6300$  matrix is taken as  $\mathbf{X}$ , and a chunk size  $c = 10000$  is used. For example, when 8 cores are used, there are  $430000 / 8 = 53750$  rows in the  $\mathbf{X}_j$  and  $\mathbf{y}_j$  on the  $j^{\text{th}}$  core, to be processed with “online” algorithms using chunk size 10000. Panels (a) and (b) respectively sketch the computation time (total = fork + join; unit: 1000 seconds) of the two reduction methods against  $m$  (the number of cores). For parallel QR reduction, while the computation time for the “fork” step (dashed line) drops with  $m$ , the time for the “join” step (dotted line) increases linearly with  $m$ , hence the gross execution time (solid line) decreases slower than the dashed line. For parallel pseudo QR reduction, the computation time for the “join” step is low and almost constant in  $m$ , thus the solid line and the dashed line almost coincide. Panels (c) and (d) sketch parallel scaling factor of both reduction methods against  $m$ , where the dashed line is the theoretical parallel scaling factor for this example (see main text for the formula), and the solid line is the realistic parallel scaling factor. The graphs confirm that parallel QR reduction has poor scaling and parallel pseudo QR reduction has good scaling.

- In “fork-join” parallel computing, QR reduction is bottlenecked by the increasingly expensive QR factorization at “join” step as more cores are used. It gives a poor practical parallel scaling. By contrast, pseudo QR reduction has very good parallel scaling.

In short, pseudo QR reduction has performance superiority in all aspects, and should be considered as the option for estimating additive models with large datasets.

### 4.3 REML estimation for additive models

Completing the discussion on model matrix reduction, I will now describe computational details for the REML estimation. For additive models,  $\mathbf{W}$  is independent of  $\hat{\beta}_{\lambda}$  (hence  $\lambda$ ), so dropping  $\log |\mathbf{W}|$  from  $\mathcal{V}_r$  gives identical estimation of  $\lambda$  and  $\phi$ . Therefore, let us consider the following REML score in subsequent derivation.

$$\mathcal{V}_r(\lambda, \phi) = (n - m) \log(2\pi\phi) + \frac{D_p(\hat{\beta}_{\lambda}, \lambda)}{\phi} + \log |\mathbf{H}_{\lambda}| - \log |\mathbf{S}_{\lambda}|_+.$$

Let  $\theta = (\log(\lambda), \log(\phi)) := (\rho, \nu)$ , minimization of  $\mathcal{V}_r$  is an unconstrained minimization w.r.t.  $\theta$ . Figure 4.9 sketches a Newton-Raphson algorithm for this minimization problem, where  $\theta^{[k]}$ ,  $\mathcal{V}_r^{[k]} = \mathcal{V}_r(\theta^{[k]})$ ,  $\nabla \mathcal{V}_r^{[k]} = \frac{\partial \mathcal{V}_r}{\partial \theta} \big|_{\theta=\theta^{[k]}}$  and  $\nabla^2 \mathcal{V}_r^{[k]} = \frac{\partial^2 \mathcal{V}_r}{\partial \theta \partial \theta'} \big|_{\theta=\theta^{[k]}}$  are respectively the parameter vector, REML score, gradient vector and Hessian matrix at the  $k^{\text{th}}$  iteration. Computation of gradient and Hessian requires the following derivatives

$$\begin{aligned} \frac{\partial \mathcal{V}_r}{\partial \rho} &= \frac{1}{\phi} \frac{\partial D_p}{\partial \rho} + \frac{\partial \log |\mathbf{H}_{\lambda}|}{\partial \rho} - \frac{\partial \log |\mathbf{S}_{\lambda}|_+}{\partial \rho}, \\ \frac{\partial^2 \mathcal{V}_r}{\partial \rho \partial \rho'} &= \frac{1}{\phi} \frac{\partial^2 D_p}{\partial \rho \partial \rho'} + \frac{\partial^2 \log |\mathbf{H}_{\lambda}|}{\partial \rho \partial \rho'} - \frac{\partial^2 \log |\mathbf{S}_{\lambda}|_+}{\partial \rho \partial \rho'}, \\ \frac{\partial \mathcal{V}_r}{\partial \nu} &= n - m - \frac{D_p}{\phi}, \quad \frac{\partial^2 \mathcal{V}_r}{\partial \rho \partial \nu} = -\frac{1}{\phi} \frac{\partial D_p}{\partial \rho}, \quad \frac{\partial^2 \mathcal{V}_r}{\partial \nu^2} = \frac{D_p}{\phi}, \end{aligned}$$

where derivatives of  $D_p(\hat{\beta}_{\lambda}, \lambda)$ ,  $\log |\mathbf{H}_{\lambda}|$  and  $\log |\mathbf{S}_{\lambda}|_+$  w.r.t.  $\rho$  are key quantities. In particular, derivatives of  $D_p(\hat{\beta}_{\lambda}, \lambda)$  demand  $\hat{\beta}_{\lambda}$  in the first place. In practice, we can arrange computations in each iteration into the following five steps:

1.  $\boxed{\log |\mathbf{S}_{\lambda}|_+}$  step computes gradient and Hessian of  $\log |\mathbf{S}_{\lambda}|_+$ , and evaluate  $\log |\mathbf{S}_{\lambda}|_+$ ;
2.  $\boxed{\hat{\beta}_{\lambda}}$  step finds  $\hat{\beta}_{\lambda} = \arg \min_{\beta} D_p(\beta, \lambda)$  and computes  $D_p(\hat{\beta}_{\lambda}, \lambda)$  and  $\log |\mathbf{H}_{\lambda}|$ ;
3.  $\boxed{\mathbf{V}_{\lambda}}$  step computes  $\mathbf{V}_{\lambda} = \mathbf{H}_{\lambda}^{-1}$ ;
4.  $\boxed{D_p(\hat{\beta}_{\lambda}, \lambda)}$  step computes gradient and Hessian of  $D_p(\hat{\beta}_{\lambda}, \lambda)$ ;
5.  $\boxed{\log |\mathbf{H}_{\lambda}|}$  step computes gradient and Hessian of  $\log |\mathbf{H}_{\lambda}|$ .

The  $\mathbf{V}_{\lambda}$  above is referred to as the *unscaled covariance matrix* of  $\hat{\beta}_{\lambda}$ , because in Bayesian point of view (see §1.1.5),  $\hat{\beta}_{\lambda}$  is an empirical Bayes estimator, whose posterior covariance matrix is  $\mathbf{V}_{\lambda}\phi$ .

REML estimation involves substantial amount of mathematics. To avoid overwhelming you, I will demonstrate how these calculations are carried out for the simple example additive model in §1.2. For ease of reference, let me first reiterate it (by copying in (1.2) and (1.11)):

$$\mathbf{y} = \mathbf{X}\beta + \beta' \mathbf{S}_{\lambda} \beta + \epsilon, \quad \epsilon \sim N(\mathbf{0}, \mathbf{W}^{-1}\phi),$$

---

```

initial value
given an initial value  $\theta^{[0]}$ , compute  $\mathcal{V}_r^{[0]}$ ,  $\nabla \mathcal{V}_r^{[0]}$ ,  $\nabla^2 \mathcal{V}_r^{[0]}$ 
Newton-Raphson algorithm
 $k = 1$ 
repeat
  Newton-type descent direction
  get a Newton step  $\Delta^{[k]} = -(\nabla^2 \mathcal{V}_r^{[k-1]})^{-1} \nabla \mathcal{V}_r^{[k-1]}$ 
   $\theta^{[k]} = \theta^{[k-1]} + \Delta^{[k]}$ , compute  $\mathcal{V}_r^{[k]}$ ,  $\nabla \mathcal{V}_r^{[k]}$ ,  $\nabla^2 \mathcal{V}_r^{[k]}$ 
  divergence test
  while  $(\mathcal{V}_r^{[k]} \geq \mathcal{V}_r^{[k-1]} + |\mathcal{V}_r^{[k-1]}| \epsilon)$ 
    step halving (backtracking line search)
     $\Delta^{[k]} = \Delta^{[k]} / 2$ 
     $\theta^{[k]} = \theta^{[k-1]} + \Delta^{[k]}$ , compute  $\mathcal{V}_r^{[k]}$ ,  $\nabla \mathcal{V}_r^{[k]}$ ,  $\nabla^2 \mathcal{V}_r^{[k]}$ 
  convergence test
  if  $(|\mathcal{V}_r^{[k]} - \mathcal{V}_r^{[k-1]}| < |\mathcal{V}_r^{[k-1]}| \epsilon)$  break
   $k += 1$ 
optimal value at convergence
 $\theta^* = \theta^{[k]}$ ,  $\mathcal{V}_r^* = \mathcal{V}_r^{[k]}$ ,  $\nabla \mathcal{V}_r^* = \nabla \mathcal{V}_r^{[k]}$ ,  $\nabla^2 \mathcal{V}_r^* = \nabla^2 \mathcal{V}_r^{[k]}$ 

```

---

**Figure 4.9:** Newton-Raphson algorithm for  $\mathcal{V}_r$  minimization.  $\epsilon$  is a small numerical tolerance for convergence and divergence test; it is a measure of relative error. A common choice is  $10^{-6}$  or  $10^{-7}$ .

where

$$\mathbf{X} = (\mathbf{X}_0 \quad \mathbf{X}_1 \quad \mathbf{X}_2), \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}, \quad \mathbf{S}_\lambda = \begin{pmatrix} \mathbf{0} & & \\ & \lambda_1 \mathbf{S}_1 & \\ & & \sum_{i=1}^{q_2} \lambda_{2i} \mathbf{S}_{2i} \end{pmatrix}.$$

In practice having more terms in an additive model does not change the methods to be presented; one just needs a loop to apply them.

During REML estimation, we will need to compute 1<sup>st</sup> and 2<sup>nd</sup> derivatives of  $\mathbf{S}_\lambda$  w.r.t.  $\rho_1 = \log(\lambda_1)$  and  $\rho_{2i} = \log(\lambda_{2i})$ . This is easy, noting that smoothing parameters enter  $\mathbf{S}_\lambda$  linearly:

$$\mathbf{S}_\lambda = \lambda_1 \begin{pmatrix} \mathbf{0} & & \\ & \mathbf{S}_1 & \\ & & \mathbf{0} \end{pmatrix} + \sum_{i=1}^{q_2} \lambda_{2i} \begin{pmatrix} \mathbf{0} & & \\ & \mathbf{0} & \\ & & \mathbf{S}_{2i} \end{pmatrix} := \exp(\rho_1) \tilde{\mathbf{S}}_1 + \sum_{i=1}^{q_2} \exp(\rho_{2i}) \tilde{\mathbf{S}}_{2i}.$$

That  $\rho$  is double indexed complicates presentation of derivatives, especially the 2<sup>nd</sup> derivatives, as I have to distinguish  $\partial \mathbf{S}_\lambda / \partial \rho_1^2$ ,  $\partial \mathbf{S}_\lambda / \partial \rho_{2i}^2$ ,  $\partial \mathbf{S}_\lambda / \partial \rho_{2j}^2$ ,  $\partial \mathbf{S}_\lambda / \partial \rho_1 \partial \rho_{2i}$ ,  $\partial \mathbf{S}_\lambda / \partial \rho_1 \partial \rho_{2j}$  and  $\partial \mathbf{S}_\lambda / \partial \rho_{2i} \partial \rho_{2j}$ . When necessary, I will index  $\rho$  and  $\tilde{\mathbf{S}}$  with a single subscript and write

$$\mathbf{S}_\lambda := \sum_l \exp(\rho_l) \tilde{\mathbf{S}}_l, \tag{4.3}$$

then the derivatives can be concisely presented by

$$\frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} = \exp(\rho_i) \tilde{\mathbf{S}}_i, \quad \frac{\partial^2 \mathbf{S}_\lambda}{\partial \rho_i \partial \rho_j} = \delta_i^j \exp(\rho_i) \tilde{\mathbf{S}}_i = \delta_i^j \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i}. \tag{4.4}$$

In the three derivative-related computation steps, double indexing style for smoothing parameters is convenient for  $\boxed{\log \|\mathbf{S}_\lambda\|_+}$  step, whereas single indexing style is convenient for  $\boxed{D_p(\hat{\boldsymbol{\beta}}_\lambda, \boldsymbol{\lambda})}$  and  $\boxed{\log \|\mathbf{H}_\lambda\|}$  steps. Readers should bear this in mind as in the forthcoming mathematical derivations for those computation steps, such indexing convention is not to be reiterated.

The rest of this section is arranged as such.

- §4.3.1 explains an initial reparametrization to  $\boldsymbol{\beta}$ , which helps reduce computational complexity of the five computation steps.



- §4.3.2 to §4.3.6 cover computational details for the five computation steps.
- §4.3.7 elaborates a special numerical issue: the stable evaluation of  $\log |\mathbf{S}_\lambda|_+$ .
- §4.3.8 summarizes FLOP count for the five computation steps and mentions the necessary post-processing to reverse the initial reparametrization when computing  $\hat{\mathbf{y}}$  or making out-of-sample prediction.

### 4.3.1 Initial reparametrization

This section demonstrates the initial parametrization to the example additive model (1.10) at the start of its REML estimation. Let  $p_i$  be the number of columns of  $\mathbf{X}_i$  (so that  $\sum_i p_i = p$ ) and  $m_i$  be the null space dimension of  $\mathbf{S}_i$  (or the number of  $\mathbf{S}_i$ 's zero eigenvalues). I will show how to reparametrize  $\beta_1$  and  $\beta_2$ , as well as how  $\mathbf{S}_\lambda$  and  $\mathbf{R}^*$  (see (4.2)) are transformed accordingly.

Consider  $\mathbf{f}_1$  first. After an eigen decomposition  $\mathbf{S}_1 = \mathbf{U}_1 \mathbf{D}_1 \mathbf{U}_1'$ , the trailing  $m_1$  eigenvalues on  $\mathbf{D}_1$ 's main diagonal are zeros. Let  $\dot{\mathbf{D}}_1$  be a modified  $\mathbf{D}_1$  where these zero eigenvalues are replaced by ones, and  $\dot{\mathbf{I}}_1$  be a modified identity matrix where the trailing  $m_1$  ones on its main diagonal are set to zeros, then the eigen decomposition can be written as  $\mathbf{S}_1 = (\mathbf{U}_1 \dot{\mathbf{D}}_1^{\frac{1}{2}}) \dot{\mathbf{I}}_1 (\mathbf{U}_1 \dot{\mathbf{D}}_1^{\frac{1}{2}})'$ . Defining a reparametrization  $\dot{\beta}_1 = \dot{\mathbf{D}}_1^{\frac{1}{2}} \mathbf{U}_1' \beta_1$ , the quadratic penalty  $\lambda_1 \beta_1' \mathbf{S}_1 \beta_1$  reduces to  $\lambda_1 \dot{\beta}_1' \dot{\mathbf{I}}_1 \dot{\beta}_1$  that resembles a ridge penalty. In fact, this reparametrization makes the link between smoothing models and mixed models more evident, as the first  $(p_1 - m_1)$  and the last  $m_1$  elements of  $\dot{\beta}_1$  can be respectively interpreted as an i.i.d. random effect (which is penalized) and some fixed effect (which is not penalized).

For  $\mathbf{f}_2$ , consider an eigen decomposition  $\sum_{i=1}^{q_2} \mathbf{S}_{2i} = \mathbf{U}_2 \mathbf{D}_2 \mathbf{U}_2'$ . As  $\mathbf{S}_2 = \sum_{i=1}^{q_2} \lambda_{2i} \mathbf{S}_{2i}$  has null space dimension  $m_2$  (for any  $\lambda_{2i}$ ), the last  $m_2$  eigenvalues in  $\mathbf{D}_2$  are zeros. Define a reparametrization  $\dot{\beta}_2 = \mathbf{U}_2' \beta_2$ , then the first  $(p_2 - m_2)$  and the last  $m_2$  elements of  $\dot{\beta}_2$  can be respectively regarded as some random effect (which is penalized) and some fixed effect (which is not penalized). Let  $\bar{\mathbf{U}}_2 = \mathbf{U}_2(1 : (p_2 - m_2))$ , then  $\dot{\mathbf{S}}_2 = \sum_{i=1}^{q_2} \lambda_{2i} \bar{\mathbf{S}}_{2i}$  is a reasonable full-rank penalty for the random effect, where  $\bar{\mathbf{S}}_{2i} = \bar{\mathbf{U}}_2' \mathbf{S}_{2i} \bar{\mathbf{U}}_2$  is a projection of  $\mathbf{S}_{2i}$  onto  $\mathbf{S}_2$ 's column space.

In above demonstration it is assumed that the null space dimension  $m_i$  of  $\mathbf{S}_i$  is known a priori. While this is true for most spline basis, it can also be numerically determined otherwise, because eigen decomposition is rank revealing. Either way, the reparametrizations above for  $\mathbf{f}_1$  and  $\mathbf{f}_2$  can be summarized by a single transformation  $\dot{\beta} = \mathbf{U} \beta$ , where

$$\mathbf{U} = \begin{pmatrix} \mathbf{I} & & \\ & \dot{\mathbf{D}}_1^{\frac{1}{2}} \mathbf{U}_1' & \\ & & \mathbf{U}_2' \end{pmatrix}, \quad \mathbf{U}^{-1} = \begin{pmatrix} \mathbf{I} & & \\ & \mathbf{U}_1 \dot{\mathbf{D}}_1^{-\frac{1}{2}} & \\ & & \mathbf{U}_2 \end{pmatrix}.$$

With the inverse transformation  $\beta = \mathbf{U}^{-1} \dot{\beta}$ , the penalized weighted residual sum of squares  $D_p(\beta, \lambda)$  (see (4.2) if you forgot what it is) can be rewritten as

$$D_p(\dot{\beta}, \lambda) = \|\mathbf{f} - \dot{\mathbf{R}}^* \dot{\beta}\|^2 + \dot{\beta}' \dot{\mathbf{S}}_\lambda \dot{\beta} + \|\mathbf{W}^{\frac{1}{2}} \mathbf{y}\|^2 - \|\mathbf{f}\|^2,$$

where  $\dot{\mathbf{R}}^* = \mathbf{R}^* \mathbf{U}^{-1}$  and

$$\dot{\mathbf{S}}_\lambda = \begin{pmatrix} \mathbf{0} & & \\ & \begin{bmatrix} \lambda_1 \mathbf{I} & \\ & \mathbf{0} \end{bmatrix} & \\ & & \begin{bmatrix} \sum_{i=1}^{q_2} \lambda_{2i} \bar{\mathbf{S}}_{2i} & \\ & \mathbf{0} \end{bmatrix} \end{pmatrix}. \quad (4.5)$$

The initial reparametrization is most intuitive in terms of the similarity transformation on  $\mathbf{S}_\lambda$ . For convenience, depending on how the penalty matrix  $\mathbf{S}_i$  of a spline  $\mathbf{f}_i$  enters  $\mathbf{S}_\lambda$ , let us call  $\mathbf{S}_i$  a *single- $\lambda$*

*block*, if it has only one smoothing parameter (like  $\mathbf{S}_1$  in the example), and a *multi- $\lambda$  block* otherwise (like  $\mathbf{S}_2$  in the example). Then in general, the initial reparametrization does the following.

- A single- $\lambda$  block is transformed to an identity-alike block, whose main diagonal elements (before multiplying  $\lambda$ ) are either one or zero. A special case is when this block is readily diagonal (for example, the ridge penalty of an i.i.d. random effect), then nothing but a scaling is needed.
- A multi- $\lambda$  block (usually rank-deficient) is projected onto a full-rank multi- $\lambda$  block with a smaller dimension, and a trailing block of zeros.

The initial reparametrization is useful in two ways.

Firstly, computing  $|\dot{\mathbf{S}}_\lambda|_+$  is more straightforward than computing  $|\mathbf{S}_\lambda|_+$ , because all non-zero diagonal blocks in  $\dot{\mathbf{S}}_\lambda$  have full rank, while similar blocks in  $\mathbf{S}_\lambda$  do not. As a result, we have

$$|\dot{\mathbf{S}}_\lambda|_+ = \lambda_1^{(p_1 - m_1)} |\bar{\mathbf{S}}_2|, \quad \log |\dot{\mathbf{S}}_\lambda|_+ = (p_1 - m_1) \rho_1 + \log |\bar{\mathbf{S}}_2|. \quad (4.6)$$

However, it is worth emphasizing that  $|\dot{\mathbf{S}}_\lambda|_+ \neq |\mathbf{S}_\lambda|_+$ , as they are respectively the log-determinant in the REML score parametrized with  $\dot{\boldsymbol{\beta}}$  and  $\boldsymbol{\beta}$ . To be precise, after the initial reparametrization, we estimate  $\dot{\boldsymbol{\beta}}$  not  $\boldsymbol{\beta}$  in Newton-Raphson iterations. We back transform  $\dot{\boldsymbol{\beta}}$  to  $\boldsymbol{\beta}$  after convergence of the iterations. I will mention this issue again in §4.3.8.

Secondly, it makes computation more efficient.  $\llbracket \log \mathbf{H}_\lambda \rrbracket$  step will need to compute  $\tilde{\mathbf{S}}_i \mathbf{V}_\lambda$  (see (4.3) if you forgot what  $\tilde{\mathbf{S}}_i$  is). Without the initial parametrization,  $\tilde{\mathbf{S}}_i$  is either of the following (double indexing style is used to distinguish a single- $\lambda$  block and a multi- $\lambda$  block):

$$\tilde{\mathbf{S}}_1 = \begin{pmatrix} \mathbf{0} & & \\ & \mathbf{S}_1 & \\ & & \mathbf{0} \end{pmatrix}, \quad \tilde{\mathbf{S}}_{2i} = \begin{pmatrix} \mathbf{0} & & \\ & \mathbf{0} & \\ & & \mathbf{S}_{2i} \end{pmatrix},$$

whereas after the initial reparametrization it is either of the following:

$$\tilde{\mathbf{S}}_1 = \begin{pmatrix} \mathbf{0} & & \\ & \begin{bmatrix} \mathbf{I} & \\ & \mathbf{0} \end{bmatrix} & \\ & & \mathbf{0} \end{pmatrix}, \quad \tilde{\mathbf{S}}_{2i} = \begin{pmatrix} \mathbf{0} & & \\ & \mathbf{0} & \\ & & \begin{bmatrix} \bar{\mathbf{S}}_{2i} & \\ & \mathbf{0} \end{bmatrix} \end{pmatrix}.$$

The matrix multiplication  $\tilde{\mathbf{S}}_1 \mathbf{V}_\lambda$  is just a row extraction on  $\mathbf{V}_\lambda$  in the latter case, but involves a dense matrix-matrix multiplication in the former.

To simplify notation, I will hereafter drop the ‘ $\cdot$ ’ from  $\dot{\boldsymbol{\beta}}$ ,  $\dot{\mathbf{R}}^*$  and  $\dot{\mathbf{S}}_\lambda$  from §4.3.2 to §4.3.7, assuming that  $\boldsymbol{\beta}$ ,  $\mathbf{R}^*$  and  $\mathbf{S}_\lambda$  have already absorbed the initial reparametrization.

### 4.3.2 Derivatives of $\log |\mathbf{S}_\lambda|_+$

Following (4.6), derivatives w.r.t.  $\rho_1$  are

$$\frac{\partial \log |\mathbf{S}_\lambda|_+}{\partial \rho_1} = (p_1 - m_1), \quad \frac{\partial^2 \log |\mathbf{S}_\lambda|_+}{\partial \rho_1^2} = 0.$$

For derivatives w.r.t.  $\boldsymbol{\rho}_2 = (\rho_{21}, \dots, \rho_{2q_2})$ , define  $\mathbf{D}_{2i} = \bar{\mathbf{S}}_2^{-1} \frac{\partial \bar{\mathbf{S}}_2}{\partial \rho_{2i}} = \exp(\rho_{2i}) \bar{\mathbf{S}}_2^{-1} \bar{\mathbf{S}}_{2i}$ , then using Harville (1997, equation (8.6), §15.8; equation (9.3), §15.9), there are

$$\begin{aligned} \frac{\partial \log |\mathbf{S}_\lambda|_+}{\partial \rho_{2i}} &= \frac{\partial \log |\bar{\mathbf{S}}_2|}{\partial \rho_{2i}} = \text{tr}(\mathbf{D}_{2i}), \\ \frac{\partial^2 \log |\mathbf{S}_\lambda|_+}{\partial \rho_{2i} \partial \rho_{2j}} &= \frac{\partial^2 \log |\bar{\mathbf{S}}_2|}{\partial \rho_{2i} \partial \rho_{2j}} = \delta_i^j \text{tr}(\mathbf{D}_{2i}) - \text{tr}(\mathbf{D}_{2i} \mathbf{D}_{2j}), \end{aligned}$$

which can be very efficiently computed using

$$\text{tr}(\mathbf{D}_{2i}) = \sum_{s=1}^p \mathbf{D}_{2i}(s, s), \quad \text{tr}(\mathbf{D}_{2i}\mathbf{D}_{2j}) = \sum_{s=1}^p \sum_{t=1}^p \mathbf{D}_{2i}(s, t)\mathbf{D}'_{2j}(s, t).$$

Furthermore, there are no 2<sup>nd</sup> derivatives “between blocks”, i.e.,

$$\frac{\partial^2 \log |\mathbf{S}_\lambda|_+}{\partial \boldsymbol{\rho}_2 \partial \rho_1} = \mathbf{0}.$$

As a result, the Hessian matrix for  $\log |\mathbf{S}_\lambda|_+$  is

$$\frac{\partial^2 \log |\mathbf{S}_\lambda|_+}{\partial \boldsymbol{\rho} \partial \boldsymbol{\rho}'} = \begin{pmatrix} 0 & \frac{\partial^2 \log |\bar{\mathbf{S}}_2|}{\partial \boldsymbol{\rho}_2 \partial \boldsymbol{\rho}_2'} \end{pmatrix}.$$

Computation of the key quantity  $\mathbf{D}_{2i}$  does not require an explicit matrix inverse  $\bar{\mathbf{S}}_2^{-1}$ . Either a QR factorization or a Cholesky factorization can be used.

- A pivoted QR factorization (see §4.1.1 if you need a revision)  $\bar{\mathbf{S}}_2 = \mathbf{Q}_2 \mathbf{R}_2^*$  yields  $\mathbf{D}_{2i} = \mathbf{R}_2^{*-1}(\mathbf{Q}_2' \bar{\mathbf{S}}_{2i})$ . Computation of the orthonormal transform can be done by a sequence of Householder transformations, as is explained in §4.1.1. This involves  $2p_2^3$  FLOP. Solving the subsequent triangular system of linear equations (see §4.1.3 if you need a revision) involves  $p_2^3$  FLOP. Since there are  $q_2$  smoothing parameters in  $\bar{\mathbf{S}}_2$ , the overall FLOP count for computing gradient and Hessian of this multi- $\lambda$  block is  $(3q_2 p_2^3 + \frac{4}{3} p_2^3)$ , where the additional  $\frac{4}{3} p_2^3$  FLOP come from the initial QR factorization.
- Since  $\bar{\mathbf{S}}_2$  is a positive-definite matrix, a more efficient method is to compute its pivoted Cholesky factorization (with pre-conditioning) (see §4.1.2 if you need a revision)  $\bar{\mathbf{S}}_2 = \mathbf{R}_2^{*'} \mathbf{R}_2^*$ , then obtain  $\mathbf{D}_{2i}$  by solving two triangular systems  $\mathbf{D}_{2i} = \mathbf{R}_2^{*-1}(\mathbf{R}_2^{*'}^{-1} \bar{\mathbf{S}}_{2i})$ . This would reduce the overall complexity to  $(2q_2 p_2^3 + \frac{1}{3} p_2^3)$  FLOP.

Given  $\mathbf{R}_2^*$ , the log-determinant  $\log |\bar{\mathbf{S}}_2|$  can be easily evaluated.

- For the QR approach with  $\mathbf{R}_2^* = \mathbf{R}_2 \mathbf{P}_2'$ , there is  $\log |\bar{\mathbf{S}}_2| = \sum_{i=1}^{p_2} \log |\mathbf{R}_2(i, i)|$ .
- For the Cholesky approach with  $\mathbf{R}_2^* = \mathbf{R}_2 \mathbf{P}_2' \mathbf{J}_2$ , there is  $\log |\bar{\mathbf{S}}_2| = 2 \sum_{i=1}^{p_2} (\log \mathbf{R}_2(i, i) + \log \mathbf{J}_2(i, i))$ .

However, the result does not always end up correct. This is a special numerical issue with a different flavour from other aspects of REML estimation. I will skip it for now and explain its stable evaluation in §4.3.7.

### 4.3.3 Estimation of $\hat{\beta}_\lambda$

Consider a pivoted Cholesky factorization (see §4.1.2 if you need a revision) for the penalty matrix  $\mathbf{S}_\lambda = \mathbf{E}_\lambda^{*'} \mathbf{E}_\lambda^*$ . Note that since  $\mathbf{S}_\lambda$  is block diagonal (see (4.5) if you forgot what it is), there is

$$\mathbf{E}_\lambda^* = \begin{pmatrix} \mathbf{0} & \begin{bmatrix} \lambda_1^{\frac{1}{2}} \mathbf{I} & \mathbf{0} \end{bmatrix} \\ \begin{bmatrix} \mathbf{E}_2^* & \mathbf{0} \end{bmatrix} \end{pmatrix}, \quad \bar{\mathbf{S}}_2 = \mathbf{E}_2^{*'} \mathbf{E}_2^*.$$

In fact,  $\mathbf{E}_\lambda^*$  will be readily available, if Cholesky method has been used for computing derivatives in the  $\lfloor \log \lfloor \mathbf{S}_\lambda \rfloor_+ \rfloor$  step. For example,  $\mathbf{E}_2^*$  here is just  $\mathbf{R}_2^*$ . Clearly there are  $(p - p_0 - m_1 - m_2)$  number of rows of zeros, which can be dropped from  $\mathbf{E}_\lambda^*$ . Now,  $D_p(\beta, \lambda)$  (see (4.2) if you forgot what it is) can be rewritten as

$$\left\| \begin{pmatrix} \mathbf{f} \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \mathbf{R}^* \\ \mathbf{E}_\lambda^* \end{pmatrix} \beta \right\|^2 + \|\mathbf{W}^{\frac{1}{2}} \mathbf{y}\|^2 - \|\mathbf{f}\|^2.$$

To find its minimizer, take a thin pivoted QR factorization (see §4.1.1 if you need a revision)

$$\begin{pmatrix} \mathbf{R}^* \\ \mathbf{E}_\lambda^* \end{pmatrix} = \mathbf{Q}_\lambda \mathbf{R}_\lambda^* = \mathbf{Q}_\lambda \mathbf{R}_\lambda \mathbf{P}_\lambda',$$

then the objective function becomes

$$\|\mathbf{f}_\lambda - \mathbf{R}_\lambda^* \beta\|^2 + \|\mathbf{W}^{\frac{1}{2}} \mathbf{y}\|^2 - \|\mathbf{f}_\lambda\|^2, \quad \mathbf{f}_\lambda = \mathbf{Q}_\lambda' \begin{pmatrix} \mathbf{f} \\ \mathbf{0} \end{pmatrix}.$$

Its minimizer of is  $\hat{\beta}_\lambda = \mathbf{R}_\lambda^{*-1} \mathbf{f}_\lambda$  which can be computed by solving a triangular system of linear equations (see §4.1.3 if you need a revision). Meanwhile,  $D_p(\hat{\beta}_\lambda, \lambda) = \|\mathbf{W}^{\frac{1}{2}} \mathbf{y}\|^2 - \|\mathbf{f}_\lambda\|^2$ . Note that the QR factorization implies a pivoted Cholesky factorization  $\mathbf{H}_\lambda = \mathbf{R}^{*'} \mathbf{R}^* + \mathbf{S}_\lambda = (\mathbf{Q}_\lambda \mathbf{R}_\lambda^*)' \mathbf{Q}_\lambda \mathbf{R}_\lambda^* = \mathbf{R}_\lambda^{*'} \mathbf{R}_\lambda^*$ , hence,  $\log |\mathbf{H}_\lambda| = 2 \sum_{i=1}^p \log |\mathbf{R}_\lambda(i, i)|$ .

Sometimes  $\mathbf{R}_\lambda$  can be singular with rank  $r < p$ . In this case, initialize  $\mathbf{z}$  as a length- $p$  vector of zeros and set  $\mathbf{z}(1:r) = \mathbf{R}_\lambda(1:r, 1:r)^{-1} \mathbf{f}_\lambda(1:r)$ , then compute  $\hat{\beta}_\lambda = \mathbf{P}_\lambda \mathbf{z}$ ,  $D_p(\hat{\beta}_\lambda, \lambda) = \|\mathbf{W}^{\frac{1}{2}} \mathbf{y}\|^2 - \|\mathbf{f}_\lambda(1:r)\|^2$  and  $\log |\mathbf{H}_\lambda| = \sum_{i=1}^r \log |\mathbf{R}_\lambda(i, i)|$ .

Computations in this step are dominated by the QR factorization of the (approximately)  $2p \times p$  augmented model matrix, which involves  $\frac{10}{3}p^3$  FLOP.

#### 4.3.4 Computing unscaled covariance matrix $\mathbf{V}_\lambda$

Given the Cholesky factorization  $\mathbf{H}_\lambda = \mathbf{R}_\lambda^{*'} \mathbf{R}_\lambda^*$ , its inverse  $\mathbf{V}_\lambda$  can be computed using methods described in §4.1.4 with  $\frac{2}{3}p^3$  FLOP. In case  $\mathbf{R}_\lambda^*$  is singular with rank  $r < p$ , initialize  $\mathbf{Z}$  as a  $p \times p$  matrix of zeros and set  $\mathbf{Z}(1:r, 1:r) = (\mathbf{R}_\lambda(1:r, 1:r)' \mathbf{R}_\lambda(1:r, 1:r))^{-1}$ , then compute  $\mathbf{V}_\lambda = \mathbf{P}_\lambda \mathbf{Z} \mathbf{P}_\lambda'$  by row and column permutations.

#### 4.3.5 Derivatives of $D_p(\hat{\beta}_\lambda, \lambda)$

As  $D_p(\hat{\beta}_\lambda, \lambda) = D(\hat{\beta}_\lambda) + \hat{\beta}_\lambda' \mathbf{S}_\lambda \hat{\beta}_\lambda$ , we consider derivatives of the two parts separately.

Using chain rule and product rule, the 1<sup>st</sup> and 2<sup>nd</sup> derivatives of  $D(\hat{\beta}_\lambda)$  w.r.t.  $\boldsymbol{\rho}$  are

$$\frac{\partial D}{\partial \rho_i} = \frac{\partial \hat{\beta}_\lambda'}{\partial \rho_i} \frac{\partial D}{\partial \hat{\beta}_\lambda}, \quad \frac{\partial^2 D}{\partial \rho_i \partial \rho_j} = \frac{\partial^2 \hat{\beta}_\lambda'}{\partial \rho_i \partial \rho_j} \frac{\partial D}{\partial \hat{\beta}_\lambda} + \frac{\partial \hat{\beta}_\lambda'}{\partial \rho_i} \frac{\partial^2 D}{\partial \hat{\beta}_\lambda \partial \hat{\beta}_\lambda'} \frac{\partial \hat{\beta}_\lambda}{\partial \rho_j}.$$

Among these quantities, derivatives of  $D(\hat{\beta}_\lambda)$  w.r.t.  $\hat{\beta}_\lambda$  can be obtained from  $D(\hat{\beta}_\lambda) = D_p(\hat{\beta}_\lambda, \lambda) - \hat{\beta}_\lambda' \mathbf{S}_\lambda \hat{\beta}_\lambda$ . Knowing (since  $\hat{\beta}_\lambda$  is the minimizer of  $D_p$ ; see (1.6) if you need a revision)

$$\frac{\partial D_p}{\partial \hat{\beta}_\lambda} = \mathbf{0}, \quad \frac{\partial^2 D_p}{\partial \hat{\beta}_\lambda \partial \hat{\beta}_\lambda'} = 2\mathbf{H}_\lambda,$$

as well as

$$\frac{\partial \hat{\beta}_\lambda' \mathbf{S}_\lambda \hat{\beta}_\lambda}{\partial \hat{\beta}_\lambda} = 2\mathbf{S}_\lambda \hat{\beta}_\lambda, \quad \frac{\partial \hat{\beta}_\lambda' \mathbf{S}_\lambda \hat{\beta}_\lambda}{\partial \hat{\beta}_\lambda \partial \hat{\beta}_\lambda'} = 2\mathbf{S}_\lambda,$$

we have

$$\frac{\partial D}{\partial \hat{\beta}_\lambda} = -2\mathbf{S}_\lambda \hat{\beta}_\lambda, \quad \frac{\partial^2 D}{\partial \hat{\beta}_\lambda \partial \hat{\beta}'_\lambda} = 2(\mathbf{H}_\lambda - \mathbf{S}_\lambda).$$

Thus,

$$\frac{\partial D}{\partial \rho_i} = -2 \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \mathbf{S}_\lambda \hat{\beta}_\lambda, \quad \frac{\partial^2 D}{\partial \rho_i \partial \rho_j} = -2 \frac{\partial^2 \hat{\beta}'_\lambda}{\partial \rho_i \partial \rho_j} \mathbf{S}_\lambda \hat{\beta}_\lambda + 2 \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} (\mathbf{H}_\lambda - \mathbf{S}_\lambda) \frac{\partial \hat{\beta}_\lambda}{\partial \rho_j}.$$

Now, using product rule, the 1<sup>st</sup> and 2<sup>nd</sup> derivatives of  $\hat{\beta}'_\lambda \mathbf{S}_\lambda \hat{\beta}_\lambda$  are

$$\begin{aligned} \frac{\partial \hat{\beta}'_\lambda \mathbf{S}_\lambda \hat{\beta}_\lambda}{\partial \rho_i} &= \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \mathbf{S}_\lambda \hat{\beta}_\lambda + \hat{\beta}'_\lambda \mathbf{S}_\lambda \frac{\partial \hat{\beta}_\lambda}{\partial \rho_i} + \hat{\beta}'_\lambda \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda = 2 \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \mathbf{S}_\lambda \hat{\beta}_\lambda + \hat{\beta}'_\lambda \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda, \\ \frac{\partial^2 \hat{\beta}'_\lambda \mathbf{S}_\lambda \hat{\beta}_\lambda}{\partial \rho_i \partial \rho_j} &= 2 \left( \frac{\partial^2 \hat{\beta}_\lambda}{\partial \rho_i \partial \rho_j} \mathbf{S}_\lambda \hat{\beta}_\lambda + \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \frac{\partial \mathbf{S}_\lambda}{\partial \rho_j} \hat{\beta}_\lambda + \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \mathbf{S}_\lambda \frac{\partial \hat{\beta}_\lambda}{\partial \rho_j} \right) + \\ &\quad \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_j} \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda + \hat{\beta}'_\lambda \frac{\partial^2 \mathbf{S}_\lambda}{\partial \rho_i \partial \rho_j} \hat{\beta}_\lambda + \hat{\beta}'_\lambda \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \frac{\partial \hat{\beta}_\lambda}{\partial \rho_j} \\ &= 2 \left( \frac{\partial^2 \hat{\beta}_\lambda}{\partial \rho_i \partial \rho_j} \mathbf{S}_\lambda \hat{\beta}_\lambda + \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \frac{\partial \mathbf{S}_\lambda}{\partial \rho_j} \hat{\beta}_\lambda + \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \mathbf{S}_\lambda \frac{\partial \hat{\beta}_\lambda}{\partial \rho_j} + \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_j} \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda \right) + \\ &\quad \hat{\beta}'_\lambda \frac{\partial^2 \mathbf{S}_\lambda}{\partial \rho_i \partial \rho_j} \hat{\beta}_\lambda. \end{aligned}$$

When combining results from the two parts, many terms cancel out and we have

$$\begin{aligned} \frac{\partial D_p}{\partial \rho_i} &= \hat{\beta}'_\lambda \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda, \\ \frac{\partial^2 D_p}{\partial \rho_i \partial \rho_j} &= 2 \left( \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \frac{\partial \mathbf{S}_\lambda}{\partial \rho_j} \hat{\beta}_\lambda + \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_i} \mathbf{H}_\lambda \frac{\partial \hat{\beta}_\lambda}{\partial \rho_j} + \frac{\partial \hat{\beta}'_\lambda}{\partial \rho_j} \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda \right) + \hat{\beta}'_\lambda \frac{\partial^2 \mathbf{S}_\lambda}{\partial \rho_i \partial \rho_j} \hat{\beta}_\lambda. \end{aligned}$$

Using (4.4), the 1<sup>st</sup> and 2<sup>nd</sup> derivatives of  $\mathbf{S}_\lambda$  w.r.t.  $\boldsymbol{\rho}$  can be substituted. It turns out that by defining  $\mathbf{A}_i = \exp(\rho_i) \tilde{\mathbf{S}}_i \hat{\beta}_\lambda$  and  $\mathbf{J}_i = \partial \hat{\beta}_\lambda / \partial \rho_i$ , the results are

$$\begin{aligned} \frac{\partial D_p}{\partial \rho_i} &= \hat{\beta}'_\lambda \mathbf{A}_i = \mathbf{A}'_i \hat{\beta}_\lambda, \\ \frac{\partial^2 D_p}{\partial \rho_i \partial \rho_j} &= 2(\mathbf{J}'_i \mathbf{A}_j + \mathbf{J}'_i \mathbf{H}_\lambda \mathbf{J}_j + \mathbf{J}'_j \mathbf{A}_i) + \delta_i^j \mathbf{A}'_i \hat{\beta}_\lambda. \end{aligned} \tag{4.7}$$

Assume that there are in total  $q$  smoothing parameters in  $\mathbf{S}_\lambda$ . If we combine all  $\mathbf{J}_i$  and  $\mathbf{A}_i$  respectively into a matrix  $\mathbf{J} = (\mathbf{J}_1 \ \mathbf{J}_2 \ \cdots \ \mathbf{J}_q)$  (this is a Jacobian matrix) and  $\mathbf{A} = (\mathbf{A}_1 \ \mathbf{A}_2 \ \cdots \ \mathbf{A}_q)$ , the derivatives can be written in a neat matrix form

$$\frac{\partial D_p}{\partial \boldsymbol{\rho}} = \mathbf{A}' \hat{\beta}_\lambda, \quad \frac{\partial^2 D_p}{\partial \boldsymbol{\rho} \partial \boldsymbol{\rho}'} = 2(\mathbf{J}' \mathbf{A} + \mathbf{J}' \mathbf{H}_\lambda \mathbf{J} + \mathbf{A}' \mathbf{J}) + \text{diag}(\mathbf{A}' \hat{\beta}_\lambda).$$

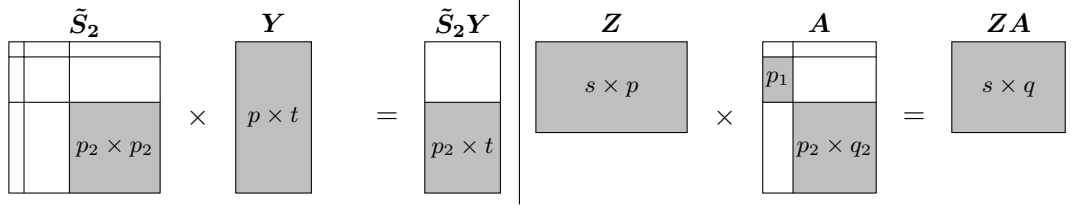
Note however that  $\mathbf{J}_i$  is unknown yet. To derive it, we apply implicit differentiation. From (1.8) we have  $\mathbf{H}_\lambda \hat{\beta}_\lambda = \mathbf{X}' \mathbf{W} \mathbf{y}$  where the RHS is independent of  $\boldsymbol{\lambda}$ . Differentiating both sides w.r.t.  $\rho_i$  gives

$$\frac{\partial \mathbf{H}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda + \mathbf{H}_\lambda \mathbf{J}_i = \mathbf{0} \quad \Rightarrow \quad \mathbf{J}_i = -\mathbf{V}_\lambda \frac{\partial \mathbf{H}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda.$$

Furthermore, (1.7) implies (as  $\mathbf{X}' \mathbf{W} \mathbf{X}$  is independent of  $\boldsymbol{\lambda}$ )

$$\frac{\partial \mathbf{H}_\lambda}{\partial \rho_i} = \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \quad \Rightarrow \quad \mathbf{J}_i = -\mathbf{V}_\lambda \frac{\partial \mathbf{S}_\lambda}{\partial \rho_i} \hat{\beta}_\lambda,$$

Using (4.4) again, we get  $\mathbf{J}_i = -\mathbf{V}_\lambda \mathbf{A}_i$ .



**Figure 4.10:** Two kinds of sparse matrix multiplications in  $[D_p(\hat{\beta}_\lambda, \lambda)]$  step and  $[\log \mathbf{H}_\lambda]$  step. In the graphs, dense blocks of a matrix are gray shaded, with its dimension printed at the centre. Unshaded blocks only contain 0. The  $p \times p$  block diagonal penalty matrix  $\tilde{\mathbf{S}}_2$  is an example of  $\tilde{\mathbf{S}}_i$  (see the end of §4.3.1 if you forgot what it is). The left panel illustrates matrix multiplication  $\tilde{\mathbf{S}}_2 \mathbf{Y}$ , where  $\mathbf{Y}$  is an arbitrary  $p \times t$  dense matrix (if  $t = 1$  then it is a vector). The resulting  $p \times t$  matrix is also sparse, with a  $p_2 \times t$  dense block. This multiplication involves  $2p_2^2 t$  FLOP. The right panel illustrates matrix multiplication  $\mathbf{Z} \mathbf{A}$ , where  $\mathbf{Z}$  is an arbitrary  $s \times p$  dense matrix, and the  $p \times q$  sparse matrix  $\mathbf{A}$  is the one that is computed in  $[D_p(\hat{\beta}_\lambda, \lambda)]$  step. The resulting  $s \times q$  matrix is dense, and the multiplication involves  $2s(p_1 + q_2 p_2)$  FLOP.

At this point we are ready to see some very surprising result. By substituting  $\mathbf{J}$  in  $\partial^2 D_p / \partial \rho \partial \rho'$  with  $\mathbf{J} = -\mathbf{V}_\lambda \mathbf{A}$ , we find that  $\mathbf{J}' \mathbf{H}_\lambda \mathbf{J} = \mathbf{A}' \mathbf{V}_\lambda \mathbf{H}_\lambda \mathbf{V}_\lambda \mathbf{A} = \mathbf{A}' \mathbf{V}_\lambda \mathbf{A}$  and  $\mathbf{J}' \mathbf{A} = \mathbf{A}' \mathbf{J} = -\mathbf{A}' \mathbf{V}_\lambda \mathbf{A}$ , so that the first two terms in the bracket cancel out! In the end, there are simply

$$\frac{\partial D_p}{\partial \rho} = \mathbf{A}' \hat{\beta}_\lambda, \quad \frac{\partial^2 D_p}{\partial \rho \partial \rho'} = 2\mathbf{J}' \mathbf{A} + \text{diag}(\mathbf{A}' \hat{\beta}_\lambda). \quad (4.8)$$

In practice, computations of  $\mathbf{A}$ ,  $\mathbf{J}$  and  $\mathbf{J}' \mathbf{A}$  can exploit sparsity and are very efficient. Figure 4.10 illustrates two kinds of matrix multiplications involving sparse matrices. In particular,  $\mathbf{A}$  is computed column by column, where computation of column  $\mathbf{A}_i$  belongs to the first kind. Computations of  $\mathbf{J}$  and  $\mathbf{J}' \mathbf{A}$  belong to the second kind. For the example additive model considered in this section, these computations respectively involve  $2q_2 p_2^2$ ,  $2p(p_1 + q_2 p_2)$  and  $2q(p_1 + q_2 p_2)$  FLOP.

#### 4.3.6 Derivatives of $\log |\mathbf{H}_\lambda|$

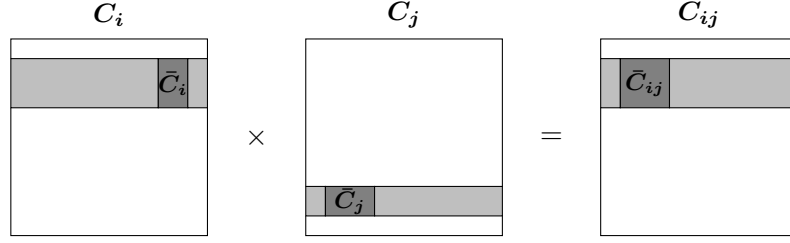
Using Harville (1997, equation (8.6), §15.8; equation (9.3), §15.9) again, it can be derived that

$$\begin{aligned} \frac{\partial \log |\mathbf{H}_\lambda|}{\partial \rho_i} &= \exp(\rho_i) \text{tr}(\mathbf{V}_\lambda \tilde{\mathbf{S}}_i), \\ \frac{\partial^2 \log |\mathbf{H}_\lambda|}{\partial \rho_i \partial \rho_j} &= \delta_i^j \exp(\rho_i) \text{tr}(\mathbf{V}_\lambda \tilde{\mathbf{S}}_i) - \exp(\rho_i + \rho_j) \text{tr}(\mathbf{V}_\lambda \tilde{\mathbf{S}}_i \mathbf{V}_\lambda \tilde{\mathbf{S}}_j), \end{aligned}$$

where the last trace also equals  $\text{tr}(\tilde{\mathbf{S}}_i \mathbf{V}_\lambda \tilde{\mathbf{S}}_j \mathbf{V}_\lambda)$ . Define  $\mathbf{C}_i = \tilde{\mathbf{S}}_i \mathbf{V}_\lambda$ , then there are

$$\begin{aligned} \frac{\partial \log |\mathbf{H}_\lambda|}{\partial \rho_i} &= \exp(\rho_i) \text{tr}(\mathbf{C}_i), \\ \frac{\partial^2 \log |\mathbf{H}_\lambda|}{\partial \rho_i \partial \rho_j} &= \delta_i^j \exp(\rho_i) \text{tr}(\mathbf{C}_i) - \exp(\rho_i + \rho_j) \text{tr}(\mathbf{C}_i \mathbf{C}_j). \end{aligned}$$

Computation of  $\mathbf{C}_i$  belongs to the first type of matrix multiplications illustrated in Figure 4.10. As a result,  $\mathbf{C}_i$  is a  $p \times p$  sparse matrix which only has one chunk of non-zero rows. Suppose these rows are indexed by a vector  $\mathbf{k}_i$ , then under the initial reparametrization, there are  $\mathbf{C}_i(\mathbf{k}_i, ) = \mathbf{V}_\lambda(\mathbf{k}_i, )$  for a single- $\lambda$  block, and  $\mathbf{C}_i(\mathbf{k}_i, ) = \tilde{\mathbf{S}}_i \mathbf{V}_\lambda(\mathbf{k}_i, )$  for a multi- $\lambda$  block. Clearly for a single- $\lambda$  block,  $\mathbf{C}_i$  needs not be computed; nor does it needs be stored as we can simply store the row index vector  $\mathbf{k}_i$ . Thus, computation of  $\mathbf{C}_i$  has a complexity that depends on the size of multi- $\lambda$  blocks only. For the example additive model considered in this section, this involves  $2q_2 p_2^2 p$  FLOP. Subsequent computations of  $\text{tr}(\mathbf{C}_i)$  and  $\text{tr}(\mathbf{C}_i \mathbf{C}_j)$  are very efficient by exploiting sparsity. For example, Figure 4.11 illustrates the computation of the latter.



**Figure 4.11:** Computation of  $\text{tr}(C_i C_j)$  in  $\lfloor \log \lfloor H_\lambda \rfloor \rfloor$  step. In the matrix product  $C_{ij} = C_i C_j$ , only the chunk of rows indexed by  $\mathbf{k}_i$  is non-zero. Hence, there is  $\text{tr}(C_{ij}) = \text{tr}(\bar{C}_{ij}) = \text{tr}(\bar{C}_i \bar{C}_j)$ , where  $\bar{C}_{ij} = C_{ij}(\mathbf{k}_i, \mathbf{k}_i)$ ,  $\bar{C}_i = C_i(\mathbf{k}_i, \mathbf{k}_j)$ , and  $\bar{C}_j = C_j(\mathbf{k}_j, \mathbf{k}_i)$ . Computation of  $\text{tr}(\bar{C}_i \bar{C}_j)$  needs no matrix multiplication. Similar to the computation of  $\text{tr}(D_{2i} D_{2j})$  in  $\lfloor \log \lfloor S_\lambda \rfloor_+ \rfloor$  step, it is just the sum of an element-wise product between  $\bar{C}_i$  and  $\bar{C}_j'$ .

#### 4.3.7 Special numerical issue: evaluation of $\log |S_\lambda|_+$

This section comes back to the evaluation of  $\log |\bar{S}_2| = \log |\sum_{i=1}^{q_2} \lambda_{2i} \bar{S}_{2i}|$  in the example additive model, or in general the log-determinant of a multi- $\lambda$  block  $M = \sum_{i=1}^q \lambda_i M_i$ , where

- $\lambda_i$ 's are positive real values;
- $M_i$ 's are positive-semidefinite matrices that do not have full rank;
- $M$  has full rank.

In §4.3.2 I have explained its basic computation methods via a pivoted QR factorization or a pivoted Cholesky factorization, with a pre-notice that it does not always give the correct answer. Now I am going to fully elaborate this matter in the following steps.

1. Present a small example where direct numerical evaluation can indeed be incorrect;
2. Explain why direct numerical evaluation goes wrong for the example;
3. Provide a way to overcome this issue and formulate this problem in general.

Let us start with the following example with  $q = 2$  penalty matrices  $M = \lambda_1 M_1 + \lambda_2 M_2$ :

$$M_1 = \begin{pmatrix} 1 & -1 & \cdot & \cdot & \cdot \\ -1 & 2 & -1 & \cdot & \cdot \\ \cdot & -1 & 1 & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}, \quad M_2 = \begin{pmatrix} \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{pmatrix},$$

where zero elements in the matrices are marked by ' $\cdot$ '. Let  $\tau = \lambda_1/\lambda_2$ , the log-determinant equals to

$$\log |\lambda_2(\tau M_1 + M_2)| = \log (|\lambda_2 I| \cdot |\tau M_1 + M_2|) = 5 \log \lambda_2 + \log |\tau M_1 + M_2|,$$

so let us focus on computing  $\log |\tau M_1 + M_2|$ . Analytical result can be easily derived:

$$|\tau M_1 + M_2| = \begin{vmatrix} \tau & -\tau & \cdot & \cdot & \cdot \\ -\tau & 2\tau & -\tau & \cdot & \cdot \\ \cdot & -\tau & \tau + 1 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{vmatrix} = \tau^2 \quad \Rightarrow \quad \log |\tau M_1 + M_2| = 2 \log \tau.$$

However, Table 4.2 shows that while numerical computation of this log-determinant is still correct as  $\tau \rightarrow 0$ , it becomes badly wrong as  $\tau \rightarrow +\infty$ . Particularly, for  $\tau = 10^{16}$  the Cholesky method returns

**Table 4.2:** Numerical computation of  $\log |\tau \mathbf{M}_1 + \mathbf{M}_2|$  can be wrong as  $\tau \rightarrow +\infty$ . See main text for  $\mathbf{M}_1$ ,  $\mathbf{M}_2$  and analytical value of the log-determinant used for this experiment.

$\tau$	1	$10^{-4}$	$10^{-8}$	$10^{-12}$	$10^{-16}$	$10^{-20}$
Analytical	0.00000	-18.42068	-36.84136	-55.26204	-73.68272	-92.10340
QR	0.00000	-18.42068	-36.84136	-55.26204	-73.68272	-92.10340
Cholesky	0.00000	-18.42068	-36.84136	-55.26204	-73.68272	-92.10340
$\tau$	1	$10^4$	$10^8$	$10^{12}$	$10^{16}$	$10^{20}$
Analytical	0.00000	18.42068	36.84136	55.26204	73.68272	92.10340
QR	0.00000	18.42068	36.84136	55.26254	75.33064	102.35677
Cholesky	0.00000	18.42068	36.84136	55.26146	-Inf	102.11145

$-\infty$ . This implies that  $\mathbf{M} = \tau \mathbf{M}_1 + \mathbf{M}_2$  becomes numerically rank-deficient, and some diagonal elements of its Cholesky factor are zeros.

To understand this issue, consider an eigen decomposition (see §4.1.5 if you need a revision)  $\mathbf{M}_1 = \mathbf{U} \mathbf{D} \mathbf{U}'$  and a similarity transformation  $\mathbf{E} = \mathbf{U}' \mathbf{M}_2 \mathbf{U}$ , then there is  $\log |\tau \mathbf{M}_1 + \mathbf{M}_2| = \log |\tau \mathbf{D} + \mathbf{E}|$ . We first show analytical results. Denote  $\dot{\mathbf{M}}_1 = \mathbf{M}_1(1:3, 1:3)$ . Its two eigenvalues  $\omega_1 = 3$  and  $\omega_2 = 1$  can be easily obtained by solving the characteristic equation  $|\dot{\mathbf{M}}_1 - \omega \mathbf{I}| = \omega(1-\omega)(\omega-3) = 0$ . It is then straightforward to obtain its eigenvectors (scaled to have unit  $L_2$ -norm)  $\boldsymbol{\nu}_1 = (\sqrt{6}/6, -\sqrt{6}/3, \sqrt{6}/6)'$  and  $\boldsymbol{\nu}_2 = (-\sqrt{2}/2, 0, \sqrt{2}/2)'$  by solving linear equations  $(\dot{\mathbf{S}}_1 - \mathbf{I})\boldsymbol{\nu}_1 = \mathbf{0}$  and  $(\dot{\mathbf{S}}_1 - 3\mathbf{I})\boldsymbol{\nu}_2 = \mathbf{0}$ . It follows that the eigen decomposition can be written as

$$\mathbf{D} = \begin{pmatrix} 3 & \cdot & \cdot & \cdot & \cdot \\ \cdot & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}, \quad \mathbf{U} = \begin{pmatrix} \sqrt{6}/6 & -\sqrt{2}/2 & \sqrt{3}/3 & \cdot & \cdot \\ -\sqrt{6}/3 & \cdot & \sqrt{3}/3 & \cdot & \cdot \\ \sqrt{6}/6 & \sqrt{2}/2 & \sqrt{3}/3 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{pmatrix},$$

where  $\mathbf{U}(:, 1) = \boldsymbol{\nu}_1$  and  $\mathbf{U}(:, 2) = \boldsymbol{\nu}_2$ . Other columns of  $\mathbf{U}$  can be determined as they must have unit  $L_2$ -norm and be orthogonal to  $\boldsymbol{\nu}_1$  and  $\boldsymbol{\nu}_2$  (so that  $\mathbf{U}$  is an orthonormal matrix). The similarity transformation is thus

$$\mathbf{E} = \mathbf{U}' \mathbf{M}_2 \mathbf{U} = \begin{pmatrix} 1/6 & \sqrt{3}/6 & \sqrt{2}/6 & \cdot & \cdot \\ \sqrt{3}/6 & 1/2 & \sqrt{6}/6 & \cdot & \cdot \\ \sqrt{2}/6 & \sqrt{6}/6 & 1/3 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{pmatrix},$$

and the determinant is

$$|\tau \mathbf{D} + \mathbf{E}| = \begin{vmatrix} 3\tau + 1/6 & \sqrt{3}/6 & \sqrt{2}/6 & \cdot & \cdot \\ \sqrt{3}/6 & \tau + 1/2 & \sqrt{6}/6 & \cdot & \cdot \\ \sqrt{2}/6 & \sqrt{6}/6 & 1/3 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{vmatrix} = \tau^2.$$

When it comes to numerical computation, the computed eigen decomposition is

$$\hat{\mathbf{D}} = \begin{pmatrix} 3 & \cdot & \cdot & \cdot & \cdot \\ \cdot & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & 2.66 \times 10^{-15} & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \end{pmatrix}, \quad \hat{\mathbf{U}} = \begin{pmatrix} 0.40825 & -0.70711 & 0.57735 & \cdot & \cdot \\ -0.81650 & \cdot & 0.57735 & \cdot & \cdot \\ 0.40825 & 0.70711 & 0.57735 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{pmatrix},$$



**Table 4.3:** Numerical computation of  $\log |\tau \mathbf{S}_1 + \mathbf{S}_2|$  becomes correct as  $\tau \rightarrow +\infty$ , if we first compute an eigen decomposition  $\mathbf{S}_1 = \hat{\mathbf{U}} \hat{\mathbf{D}} \hat{\mathbf{U}}'$ , detect the numerical rank of  $\mathbf{S}_1$ , set the spurious eigenvalues in  $\hat{\mathbf{D}}$  to zeros, then compute  $\log |\tau \hat{\mathbf{D}} + \hat{\mathbf{U}}' \mathbf{S}_2 \hat{\mathbf{U}}|$  following a similarity transformation.

	$\tau$	1	$10^4$	$10^8$	$10^{12}$	$10^{16}$	$10^{20}$
Analytical	0.00000	18.42068	36.84136	55.26204	73.68272	92.10340	
QR	0.00000	18.42068	36.84136	55.26204	73.68272	92.10340	
Cholesky	0.00000	18.42068	36.84136	55.26204	73.68272	92.10340	

and the computed similarity transformation is

$$\hat{\mathbf{E}} = \begin{pmatrix} 0.16667 & 0.28868 & 0.23570 & \cdot & \cdot \\ 0.28868 & 0.50000 & 0.40825 & \cdot & \cdot \\ 0.23570 & 0.40825 & 0.33333 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{pmatrix}.$$

We see that although  $\hat{\mathbf{U}}$  and  $\hat{\mathbf{E}}$  are identical to  $\mathbf{U}$  and  $\mathbf{E}$ , the third computed eigenvalue, denoted by  $\delta$ , is not strictly zero, but a small value at the magnitude of the machine precision (about  $2.22 \times 10^{-16}$ ). As a result, we get

$$|\tau \hat{\mathbf{D}} + \hat{\mathbf{E}}| = \begin{vmatrix} 3\tau + 1/6 & \sqrt{3}/6 & \sqrt{2}/6 & \cdot & \cdot \\ \sqrt{3}/6 & \tau + 1/2 & \sqrt{6}/6 & \cdot & \cdot \\ \sqrt{2}/6 & \sqrt{6}/6 & \delta\tau + 1/3 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot & 1 \end{vmatrix} = \tau^2 + \delta\tau(3\tau^2 + \frac{5}{3}\tau + \frac{1}{36}),$$

which deviates from the true quantity by  $\delta\tau(3\tau^2 + \frac{5}{3}\tau + \frac{1}{36})$ . As  $\tau \rightarrow +\infty$ , this error is toward  $+\infty$ .

The above demonstration also hints a way to suppress such undesired numerical property. If we, after the computed eigen decomposition, manually set  $\delta = 0$ , then  $|\tau \hat{\mathbf{D}} + \hat{\mathbf{E}}|$  would give the correct result. Table 4.3 verifies this.

Suppose that  $\mathbf{S}_1$  has rank  $r < t$ , then  $\tau \hat{\mathbf{D}} + \hat{\mathbf{E}}$  can be partitioned as

$$\begin{pmatrix} \tau \hat{\mathbf{D}}_{11} + \hat{\mathbf{E}}_{11} & \hat{\mathbf{E}}_{12} \\ \hat{\mathbf{E}}_{21} & \tau \hat{\mathbf{D}}_{22} + \hat{\mathbf{E}}_{22} \end{pmatrix},$$

where  $\hat{\mathbf{D}}_{11} = \hat{\mathbf{D}}(1 : r, 1 : r)$  contains the correct eigenvalues of  $\mathbf{S}_1$  while  $\hat{\mathbf{D}}_{22}$  contains spurious eigenvalues with  $|\hat{\mathbf{D}}_{22}(i, i)| < \hat{\mathbf{D}}_{11}(1, 1)\epsilon$ . Since  $\tau \hat{\mathbf{D}}_{11} + \hat{\mathbf{E}}_{11}$  has full rank, then using the general result

$$\begin{vmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{vmatrix} = |\mathbf{A}_{11}| \cdot |\mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12}|,$$

there is

$$\log |\tau \hat{\mathbf{D}} + \hat{\mathbf{E}}| = \log |\tau \hat{\mathbf{D}}_{11} + \hat{\mathbf{E}}_{11}| + \log |\tau \hat{\mathbf{D}}_{22} + \hat{\mathbf{E}}_{22} - \hat{\mathbf{E}}_{21}(\tau \hat{\mathbf{D}}_{11} + \hat{\mathbf{E}}_{11})^{-1} \hat{\mathbf{E}}_{12}|.$$

As  $\tau \rightarrow +\infty$ , the last log-determinant would approach  $\log |\tau \hat{\mathbf{D}}_{22} + \hat{\mathbf{E}}_{22}|$ , which would only be correct if we set  $\hat{\mathbf{D}}_{22} = \mathbf{0}$ , which we call a “truncation”. Otherwise, an upper bound of the wrong value is

$$\log |\tau \hat{\mathbf{D}}_{22} + \hat{\mathbf{E}}_{22}| < \log |\tau \hat{\mathbf{D}}_{11}(1, 1)\epsilon \mathbf{I} + \hat{\mathbf{E}}_{22}| = \sum_i \log (\tau \hat{\mathbf{D}}_{11}(1, 1)\epsilon + e_i),$$

where  $e_i$  is  $\hat{\mathbf{E}}_{22}$ 's eigenvalue, which is also  $\hat{\mathbf{E}}$ 's and hence  $\mathbf{S}_2$ 's eigenvalue. An upper bound for the error is then  $\sum_i \log (\tau \hat{\mathbf{D}}_{11}(1, 1)\epsilon / e_i + 1)$ . Obviously, whenever the maximum eigenvalue of  $\mathbf{S}_1$  times  $\tau\epsilon$  is much larger than the smallest eigenvalue of  $\mathbf{S}_2$ , a situation we call a “bad scaling” between  $\mathbf{S}_1$  and  $\mathbf{S}_2$ , the error bound is non-negligible. This may give some heuristic on when a similarity transformation and a “truncation” are possibly needed for stable evaluation of a log-determinant. A more-than-sufficient condition is as such. Compute Frobenius norm (an upper bound for maximum

**Table 4.4:** FLOP count for the five computation steps in computing  $\nabla\mathcal{V}_r$  and  $\nabla^2\mathcal{V}_r$ . See main text for meanings of  $l$ ,  $p_i$ ,  $q_i$ ,  $p$  and  $q$ .

$\lfloor \log  \mathbf{S}_\lambda _+ \rfloor$	$\lfloor \hat{\beta}_\lambda \rfloor$	$\lfloor \mathbf{V}_\lambda \rfloor$	$\lfloor D_p(\hat{\beta}_\lambda, \lambda) \rfloor$	$\lfloor \log  \mathbf{H}_\lambda  \rfloor$
$\sum_{\{i: q_i > 1\}} (3q_i + \frac{4}{3})p_i^3$	$\frac{10}{3}p^3$	$\frac{2}{3}p^3$	$\sum_{\{i: q_i > 1\}} 2q_i p_i^2 + 2(p+q) \sum_{i=1}^l q_i p_i$	$\sum_{\{i: q_i > 1\}} 2q_i p_i^2 p$

eigenvalue)  $\|\mathbf{S}_1\|_F$  and  $\|\mathbf{S}_2\|_F$  of  $\mathbf{S}_1$  and  $\mathbf{S}_2$ . If  $\lambda_1\|\mathbf{S}_1\|_F/\lambda_2\|\mathbf{S}_2\|_F > \epsilon$ , then perform similarity transformation and “truncation”.

If there are  $q > 2$  penalty matrices in a multi- $\lambda$  block, the above idea for  $q = 2$  case can be generalized. We refer readers to Wood (2011, Appendix B) for such an algorithm. In the end, there exists a single similarity transformation  $\hat{\mathbf{S}} = \hat{\mathbf{U}}'\mathbf{S}\hat{\mathbf{U}}$ , and evaluation of  $\log |\hat{\mathbf{S}}|$  is numerically stable. Such transformation also guarantees that  $\hat{\mathbf{S}}^{-1}$  is computable for calculations of derivatives. Note that a transformation to each individual penalty matrix  $\hat{\mathbf{S}}_i = \hat{\mathbf{U}}'\mathbf{S}_i\hat{\mathbf{U}}$  is needed. In addition, in the subsequent  $\lfloor \hat{\beta}_\lambda \rfloor$  step,  $\mathbf{R}^*$  needs to be transformed accordingly.

#### 4.3.8 Summary

In §4.3.2 to §4.3.6 I have given FLOP count of the five computation steps in computing  $\nabla\mathcal{V}_r$  and  $\nabla^2\mathcal{V}_r$  for the example additive model in §4.3.1. Now I will generalize these results. Suppose we have an additive model with  $\mathbf{X}\boldsymbol{\beta} = \mathbf{X}_0\boldsymbol{\beta}_0 + \sum_{i=1}^l \mathbf{f}_i$ , where there are some parametric term  $\boldsymbol{\beta}_0$  and  $l$  splines. For the  $i^{\text{th}}$  spline  $\mathbf{f}_i = \mathbf{X}_i\boldsymbol{\beta}_i$ , assume it has  $p_i$  coefficients in  $\boldsymbol{\beta}_i$  and  $q_i$  smoothing parameters in its penalty matrix  $\mathbf{S}_i$ . Note that  $q_i$  may be any non-negative integers. If  $q_i = 0$ , then  $\mathbf{f}_i$  is a pure regression spline without penalization. If  $q_i = 1$ , then  $\mathbf{S}_i$  is a single- $\lambda$  block.  $q_i > 1$  implies that  $\mathbf{S}_i$  is a multi- $\lambda$  block. The total number of coefficients in  $\boldsymbol{\beta}$  is  $p = \sum_{i=0}^l p_i$  (don’t forget the parametric term) and the total number of smoothing parameters in  $\mathbf{S}_\lambda$  is  $q = \sum_{i=1}^l q_i$ . For this general additive model, computational complexity of all steps is summarized in Table 4.4. Results for  $\lfloor \log |\mathbf{S}_\lambda|_+ \rfloor$  and  $\lfloor \log |\mathbf{H}_\lambda| \rfloor$  steps only depend on the size of multi- $\lambda$  blocks in  $\mathbf{S}_\lambda$ , as the summation is over  $\{i : q_i > 1\}$ . In addition, the result for  $\lfloor \log |\mathbf{S}_\lambda|_+ \rfloor$  step is for QR factorization approach; for Cholesky factorization approach, the sequence to be summed is  $(2q_i + \frac{1}{3})p_i^3$ . In general,  $\lfloor \hat{\beta}_\lambda \rfloor$  step is the most expensive step, and its  $O(p^3)$  complexity is often reported as an approximate complexity for REML estimation, in contrast to the  $O(np^2)$  complexity for model matrix reduction.

Apart from parameter estimation, prediction of response variable  $\mathbf{y}$ , i.e., computation of  $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}_\lambda$ , is also an important part of inference. This is a straightforward matrix-vector multiplication, but to avoid forming the whole  $\mathbf{X}$ , the row chunking strategy used in model matrix reduction (see §4.2 if you need a revision) can be likewise applied, with further possibility of parallel computing. In addition, prediction requires the original parametrization of the coefficient vector. However, the initial reparametrization before REML estimation means that at convergence of the Newton-Raphson algorithm, estimated coefficients and unscaled covariance matrix are actually associated with  $\tilde{\beta}$  not  $\beta$ . Denote these estimates by  $\tilde{\beta}_\lambda$  and  $\tilde{\mathbf{V}}_\lambda$ , back transformations  $\hat{\beta}_\lambda = \mathbf{U}^{-1}\tilde{\beta}_\lambda$  and  $\mathbf{V}_\lambda = \mathbf{U}^{-1}\tilde{\mathbf{V}}_\lambda\mathbf{U}^{-1}$  are a necessary post-processing. In case a secondary reparametrization for stable evaluation of  $\log |\mathbf{S}_\lambda|_+$  also exists, it needs to be reversed in a similar way and should be done in the first place.

## 4.4 Estimation of GAMs via “performance iteration”

A GAM can be expressed as a penalized generalized linear model (GLM), or a generalized linear mixed model (GLMM), as follows:

$$\mathbf{y}|\boldsymbol{\beta} \sim \text{EF}(\boldsymbol{\mu}, \boldsymbol{\phi}), \quad \boldsymbol{\eta} = g(\boldsymbol{\mu}) = \mathbf{X}\boldsymbol{\beta}, \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \mathbf{S}_\lambda^{-1}\boldsymbol{\phi}),$$

---

Initialize  $k = 0$ ,  $D_p^{[0]} = 0$ ,  $\boldsymbol{\mu}^{[0]} = \mathbf{y} + \boldsymbol{\xi}$  and  $\boldsymbol{\eta}^{[0]} = g(\boldsymbol{\mu}^{[0]})$ , then

1. obtain pseudo data vector  $\mathbf{z}^{[k]}$ , with elements  $z_i^{[k]} = g'(\mu_i^{[k]})(y_i - \mu_i^{[k]}) + \eta_i^{[k]}$ , and a diagonal weight matrix  $\mathbf{W}^{[k]}$ , with diagonal elements  $w_{ii}^{[k]} = V(\mu_i^{[k]})^{-1}g'(\mu_i^{[k]})^{-2}$ ;
  2. solve a penalized weighted least squares problem  $\boldsymbol{\beta}^{[k+1]} = \arg\max_{\boldsymbol{\beta}} D_p(\boldsymbol{\beta}, \boldsymbol{\lambda})$ , where  $D_p(\boldsymbol{\beta}, \boldsymbol{\lambda}) = \|(\mathbf{W}^{[k]})^{\frac{1}{2}}\mathbf{z}^{[k]} - (\mathbf{W}^{[k]})^{\frac{1}{2}}\mathbf{X}\boldsymbol{\beta}\|^2 + \boldsymbol{\beta}'\mathbf{S}_{\boldsymbol{\lambda}}\boldsymbol{\beta}$ ;
  3. evaluate  $D_p^{[k+1]} = D_p(\boldsymbol{\beta}^{[k+1]}, \boldsymbol{\lambda})$ ,  $\boldsymbol{\eta}^{[k+1]} = \mathbf{X}\boldsymbol{\beta}^{[k+1]}$  and  $\boldsymbol{\mu}^{[k+1]} = g^{-1}(\boldsymbol{\eta}^{[k+1]})$ ;
  4. if  $|D_p^{[k+1]} - D_p^{[k]}| < (0.1 + |D_p^{[k]}|)\epsilon$ , then iteration has converged and go to 5; otherwise  $k += 1$  and go to 1.
  5. set  $\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}} = \boldsymbol{\beta}^{[k+1]}$ ,  $\hat{\boldsymbol{\eta}}_{\boldsymbol{\lambda}} = \boldsymbol{\eta}^{[k+1]}$ ,  $\hat{\boldsymbol{\mu}}_{\boldsymbol{\lambda}} = \boldsymbol{\mu}^{[k+1]}$ ,  $\mathbf{W}_{\boldsymbol{\lambda}} = \mathbf{W}^{[k]}$ ,  $\mathbf{z}_{\boldsymbol{\lambda}} = \mathbf{z}^{[k]}$  and  $D_p(\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}, \boldsymbol{\lambda}) = \|\mathbf{W}_{\boldsymbol{\lambda}}^{\frac{1}{2}}\mathbf{z}_{\boldsymbol{\lambda}} - \mathbf{W}_{\boldsymbol{\lambda}}^{\frac{1}{2}}\mathbf{X}\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}\|^2 + \hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}'\mathbf{S}_{\boldsymbol{\lambda}}\hat{\boldsymbol{\beta}}_{\boldsymbol{\lambda}}$ .
- 

**Figure 4.12:** P-IRLS algorithm for GAM estimation with known  $\boldsymbol{\lambda}$ .  $\epsilon$  is some small numerical tolerance.  $\boldsymbol{\xi}$  is a vector of small quantities (often zeros) to ensure that  $g(\boldsymbol{\mu}^{[0]})$  exists.  $g'(\mu)$  is the 1<sup>st</sup> derivative of the link function  $g(\mu)$ .  $V(\mu)$  is the variance function for the exponential family distribution. Note that at convergence (step 5) of the algorithm, all quantities have dependency on  $\boldsymbol{\lambda}$ .

where EF is an exponential family distribution with known or unknown scale parameter  $\phi$ , and  $g$  is a known smooth monotonic link function that transforms  $\boldsymbol{\mu} = \mathbf{E}(\mathbf{y}|\boldsymbol{\beta})$  to the linear predictor  $\boldsymbol{\eta}$ . If all smoothing parameters in  $\boldsymbol{\lambda}$  are known, then  $\boldsymbol{\beta}$  can be estimated by maximizing the penalized log-likelihood (which in Bayesian point of view, differs from the log-posterior by some  $\boldsymbol{\lambda}$ -related constant):

$$l_p(\boldsymbol{\beta}, \boldsymbol{\lambda}) = l(\boldsymbol{\beta}) - \frac{\boldsymbol{\beta}'\mathbf{S}_{\boldsymbol{\lambda}}\boldsymbol{\beta}}{2\phi}.$$

Numerical maximization often adopts a *penalized iteratively re-weighted least squares (P-IRLS)* algorithm, as presented in Figure 4.12.

When  $\boldsymbol{\lambda}$  is unknown and also needs to be estimated, step 2 of the P-IRLS algorithm may be replaced by estimation of an additive model

$$\mathbf{z}^{[k]}|\boldsymbol{\beta} \sim N(\mathbf{X}\boldsymbol{\beta}, (\mathbf{W}^{[k]})^{-1}\phi), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \mathbf{S}_{\boldsymbol{\lambda}}^{-1}\phi),$$

using methods explained in previous sections of this Chapter. The resulting algorithm resembles the *penalized quasi-likelihood (PQL)* estimation (Breslow and Clayton, 1993; Tuerlinckx et al., 2010) of a GLMM. It may also be seen as a variant of the “performance iteration” proposed by Gu (1992) for GAM estimation using generalized cross-validation. The normality assumption for  $\mathbf{z}^{[k]}$  is clearly false. However, after the QR reduction  $(\mathbf{W}^{[k]})^{\frac{1}{2}}\mathbf{z}^{[k]} = \mathbf{Q}^{[k]'}\mathbf{R}^{*[k]}$ ,  $D_p(\boldsymbol{\beta}, \boldsymbol{\lambda})$  in step 2 of the P-IRLS algorithm can be rewritten into a form similar to (4.2), but now with  $\mathbf{f}^{[k]} = \mathbf{Q}^{[k]}(\mathbf{W}^{[k]})^{\frac{1}{2}}\mathbf{z}^{[k]}$ . This implies that the additive model for  $\mathbf{z}^{[k]}$  may be alternatively written as the following one for  $\mathbf{f}^{[k]}$ :

$$\mathbf{f}^{[k]}|\boldsymbol{\beta} \sim N(\mathbf{R}^{*[k]}\boldsymbol{\beta}, \mathbf{I}\phi), \quad \boldsymbol{\beta} \sim N(\mathbf{0}, \mathbf{S}_{\boldsymbol{\lambda}}^{-1}\phi),$$

For large datasets there is  $n \gg p$ , so asymptotic normality for  $\mathbf{f}^{[k]}$  is justified under central limit theorem.

“Performance” iteration is not the only method for GAM estimation. Alternatively, we may attempt a number of candidate smoothing parameters  $\boldsymbol{\lambda}^{[1]}, \boldsymbol{\lambda}^{[2]}, \dots$ , and for each trial the GAM is estimated with the P-IRLS algorithm. Then we choose the best  $\boldsymbol{\lambda}^{[i]}$ , which minimizes the REML score  $\mathcal{V}_r(\boldsymbol{\lambda}, \phi) = -2\log \int P(\mathbf{y}|\boldsymbol{\beta})P(\boldsymbol{\beta})d\boldsymbol{\beta}$ . This crude “grid search” can be replaced by numerical minimization via a Newton-Raphson algorithm, if  $\nabla\mathcal{V}_r$  and  $\nabla^2\mathcal{V}_r$  can be derived. The resulting estimation method is termed the “outer iteration” in `mgcv`, where we have an outer loop for  $\boldsymbol{\lambda}$  and an inner loop that iterates the P-IRLS algorithm till convergence for each trial  $\boldsymbol{\lambda}$ . For non-Gaussian distribution, the integral in  $\mathcal{V}_r$  is not in a closed form. We may approximate the integrand as follows, so that the

**Table 4.5:** “outer iteration” and “performance iteration” for GAM estimation

“outer iteration”	$\left( \begin{array}{c} \lambda^{[i]} \\ \left( \begin{array}{c} \mathbf{W}^{[k]}, \mathbf{z}^{[k]}, \beta^{[k]}, D_p^{[k]} \end{array} \right) \\ \mathcal{V}_r^{[i]}, \nabla \mathcal{V}_r^{[i]}, \nabla^2 \mathcal{V}_r^{[i]} \end{array} \right)$
“performance iteration”	$\left( \begin{array}{c} \mathbf{W}^{[k]}, \mathbf{z}^{[k]} \\ \left( \lambda^{[i]}, \beta^{[i]}, D_p^{[i]}, \mathcal{V}_r^{[i]}, \nabla \mathcal{V}_r^{[i]}, \nabla^2 \mathcal{V}_r^{[i]} \right) \end{array} \right)$

integral of the approximated integrand has a closed form.

$$\begin{aligned}
P(\mathbf{y}|\beta)P(\beta) &= \exp\{l(\beta)\}(2\pi\phi)^{-\frac{p-m}{2}}|\mathbf{S}_\lambda|_+^{\frac{1}{2}}\exp\{-\frac{\beta'\mathbf{S}_\lambda\beta}{2\phi}\} \\
&= (2\pi\phi)^{-\frac{p-m}{2}}|\mathbf{S}_\lambda|_+^{\frac{1}{2}}\exp\{l_p(\beta, \lambda)\} \\
&= (2\pi\phi)^{-\frac{p-m}{2}}|\mathbf{S}_\lambda|_+^{\frac{1}{2}}\exp\{l_s(\phi) - \frac{D_p(\beta, \lambda)}{2\phi}\} \\
&\approx (2\pi\phi)^{-\frac{p-m}{2}}|\mathbf{S}_\lambda|_+^{\frac{1}{2}}\exp\{l_s(\phi) - \frac{D_p(\hat{\beta}_\lambda, \lambda) + (\beta - \hat{\beta}_\lambda)'\mathbf{H}_\lambda(\beta - \hat{\beta}_\lambda)}{2\phi}\}.
\end{aligned}$$

From line 2 to line 3, we substitute  $l_p(\hat{\beta}_\lambda, \lambda)$  for  $D_p(\hat{\beta}_\lambda, \lambda)$  using  $D_p(\hat{\beta}_\lambda, \lambda) = 2(l_s(\phi) - l_p(\hat{\beta}_\lambda, \lambda))\phi$ , where  $l_s(\phi)$  is the *saturated log-likelihood*. From line 3 to line 4, we apply a Taylor expansion of  $D_p(\beta, \lambda)$  at its minimum point  $\hat{\beta}_\lambda$ . Integrating the approximated integrand gives a Laplace approximation to the integral:

$$\int P(\mathbf{y}|\beta)P(\beta)d\beta = |\mathbf{S}_\lambda|_+^{\frac{1}{2}}\exp\{l_s(\phi) - \frac{D_p(\hat{\beta}_\lambda, \lambda)}{2\phi}\}(2\pi\phi)^{\frac{m}{2}}|\mathbf{H}_\lambda|^{-\frac{1}{2}},$$

hence the REML score is

$$\mathcal{V}_r(\lambda, \phi) = -m\log(2\pi\phi) - 2l_s(\phi) + \frac{D_p(\hat{\beta}_\lambda, \lambda)}{\phi} + \log|\mathbf{H}_\lambda| - \log|\mathbf{S}_\lambda|_+.$$

This expression looks like the REML score (1.9) for additive models, but is in fact substantially different. Notably, at convergence of P-IRLS, there is  $\mathbf{H}_\lambda = \mathbf{X}'\mathbf{W}_\lambda\mathbf{X} + \mathbf{S}_\lambda$  where the weight matrix has dependency on  $\lambda$ . Furthermore, the dependency of  $D_p(\hat{\beta}_\lambda, \lambda)$  on  $\lambda$  is not just via its dependency on  $\hat{\beta}_\lambda$ , but also via its dependency on  $\mathbf{z}_\lambda$  and  $\mathbf{W}_\lambda$ . As a result, derivatives of  $D_p(\hat{\beta}_\lambda, \lambda)$  and  $\log|\mathbf{H}_\lambda|$  w.r.t.  $\rho = \log(\lambda)$  would be computed in a different way than what is presented in §4.3.5 and §4.3.6 for additive models. See Wood (2011) for full details.

Table 4.5 shows an intuitive comparison between the structures of “outer iteration” and “performance iteration”. Coverage of “outer iteration” is guaranteed as a minimum is well defined for both of its outer loop and inner loop. However, since the inner loop needs to update the weighted least squares problem, its computational complexity is  $O(np^2) + O(p^3)$ . For large  $n$  (and possibly big  $p$  as well), “outer” iteration is expensive or even infeasible. “Performance iteration” is in principle less costly than “outer iteration”, as its inner loop only has a  $O(p^3)$  complexity once the QR reduction is done in the outer loop. But it does not always converge, since the REML score is associated with a working additive model that changes at each iteration of the outer loop. It is sometimes possible for the algorithm to cycle around a small set of smoothing parameter, coefficient combinations without ever converging. In this situation, it is helpful to reduce  $k$  values of splines to reduce the flexibility of the model. Still, “performance iteration” is useful for GAM estimation with large datasets.

## 4.5 Summary

Note that it should be emphasized that that  $\mathbf{W}$  is loop-invariant is the premise for applying model matrix reduction. This is generally true for non-weighted additive models, or models with known weights. If for example, the model error is an AR(1) process, the autocorrelation coefficient has to be known. In other words, if this autocorrelation coefficient needs be estimated, it should not be estimated jointly with smoothing parameters in the REML estimation loop, otherwise the condition for applying model matrix reduction is invalidated.

## Chapter 5

# Optimizing GAM computations for high performance computing

In the last Chapter I reviewed the computational engine for estimating an additive model. The aim of this Chapter is to investigate where it can be potentially optimized, so that the computational hurdle encountered in Table 3.11 for fitting logBS models can be overcome. To start with, I break down the model fitting time in that table to reveal how much time is spent for pseudo QR reduction and each step of REML estimation. See Table 5.1. The rest of this Chapter is about how to, step by step, obtain better figures in Table 5.2. Here is a brief outline of these sections.

- §5.1 points out a few inefficiencies in current `bam`'s implementation for REML estimation. As mentioned at the beginning of §4, I realized a few places where the design idea could be better implemented when reviewing the computational engine. So in this section I will benchmark my improved version with the original version. As you can see from Table 5.2, the execution time for  $\log[\mathbf{S}_{\lambda}]_{+}$  and  $D_p(\hat{\beta}_{\lambda}, \lambda)$  steps are successfully reduced.
- §5.2 introduces a scientific library called *Basic Linear Algebra Subroutines (BLAS)* and its optimized distributions, explains concepts of block algorithms and data caching that are related to high performance computing, then exploits such library for GAM computations. As you can see from Table 5.2, except for the  $\mathbf{V}_{\lambda}$  step, this yields big speedup.
- §5.3 identifies why BLAS has failed to boost  $\mathbf{V}_{\lambda}$  step. It turns out that the underlying C code in `bam` for  $\mathbf{V}_{\lambda}$  computation is not using BLAS as building blocks so changing BLAS has no performance impact on it. I then re-implement this step using an existing LAPACK routine, and as you can see from Table 5.2, this delivers another big speedup.
- §5.4 explains how  $\hat{\beta}_{\lambda}$  step, the slowest part after previous optimizations, can be accelerated. The key is to solve penalized least squares with pivoted Cholesky factorization instead of pivoted QR factorization, which not only involves less FLOP count, but also better benefits from the performance gains from using an optimized BLAS.
- §5.5 makes a final push on computational speed by using parallel computing. Different parallel computing strategies are used for pseudo QR reduction and REML estimation. For the former, the “fork-join” method previously explained in §4.2.3 is used, and for the latter, multi-threading option in an optimized BLAS is enabled. Table 5.2 demonstrates the computation time for model 3.6 when 8 CPU cores are used. Model fitting time is now down to just about 6 minutes.

In §5.5, it will also be demonstrated that computational methods developed in this Chapter will be sufficiently fast to estimate a joint model for logBS from all days of week in 1967. To demonstrate the computational superiority of such new GAM fitting methods, a quick comparison with INLA is made in the summary section. It turns out that fitting daily logBS models via INLA is completely infeasible.

**Table 5.1:** A break-down of model fitting time (in seconds) previously reported in Table 3.11 for model 3.3b, model 3.5 and model 3.6. Two remarks: 1) Pseudo QR reduction is used as model matrix reduction method, as I have concluded in §4.2.4 that this method outperforms QR reduction; 2) The  $\hat{\mathbf{y}}$  column refers to the computation of fitted values  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda}$  and residuals of a model, as well as other simple statistics like estimated residual variance, AIC, etc. These models are fitted on an Intel Xeon E5-2650 v2 workstation (see Appendix A for hardware information), with R software linked to the reference BLAS. BLAS will be shortly introduced in §5.2. The early reference to it here is just to provide full information on how those timing statistics were collected.

	<i>Manchester 11</i>	annual mean logBS	Monday logBS in 1967
	model 3.3b $n = 15707$ $p = 732$	model 3.5 $n = 24239$ $p = 6750$	model 3.6 $n = 54386$ $p = 10932$
pseudo QR reduction	96.66	704.91	4046.42
$\log \ \mathbf{S}_{\lambda}\ _+$	95.76	2518.72	7732.10
$\ \hat{\boldsymbol{\beta}}_{\lambda}\ $	174.72	4232.64	23509.83
$\ \mathbf{V}_{\lambda}\ $	28.56	842.16	4356.09
$D_p(\hat{\boldsymbol{\beta}}_{\lambda}, \boldsymbol{\lambda})$	16.80	70.24	186.91
$\log \ \mathbf{H}_{\lambda}\ $	31.92	1523.20	6043.37
$\hat{\mathbf{y}}$	0.73	35.81	87.88
total	442.15	9927.68	45962.61

**Table 5.2:** A summary of how modelling estimation is accelerated by using method described in each section of this Chapter. Numbers in the table are execution time measured in seconds. A blank field in the table means that the execution time is not speeded up (nor slowed down of course) by a method, so it is as same as the last non-empty value in a column. All models are fitted on an Intel Xeon E5-2650 v2 workstation (see Appendix A for hardware information). Starting from §5.2, OpenBLAS is linked to R.

(a) Results for model 3.3b for logBS time series from *Manchester 11*. This model has an AR(1) component to estimate, so there is an outer iteration via golden-section search.

	golden-section search converges in 14 steps							
	reduction	REML estimation converges in 12 Newton steps						
		$\log \ \mathbf{S}_{\lambda}\ _+$	$\ \hat{\boldsymbol{\beta}}_{\lambda}\ $	$\ \mathbf{V}_{\lambda}\ $	$D_p(\hat{\boldsymbol{\beta}}_{\lambda}, \boldsymbol{\lambda})$	$\log \ \mathbf{H}_{\lambda}\ $	$\hat{\mathbf{y}}$	total
initial	93.66	95.76	174.72	28.56	16.80	31.92	0.73	442.15
§5.1		16.80			0.60			346.99
§5.2	14.56	2.52	45.96		0.44	4.37	0.37	96.78
§5.3				4.70				72.92
§5.4			4.87					31.83

(b) Results for model 3.5 for annual mean logBS. This model has no AR(1) component to estimate.

	REML estimation converges in 8 Newton steps						$\hat{\mathbf{y}}$	total
	reduction	$\log \ \mathbf{S}_{\lambda}\ _+$	$\ \hat{\boldsymbol{\beta}}_{\lambda}\ $	$\ \mathbf{V}_{\lambda}\ $	$D_p(\hat{\boldsymbol{\beta}}_{\lambda}, \boldsymbol{\lambda})$	$\log \ \mathbf{H}_{\lambda}\ $		
initial	704.91	2518.72	4232.64	842.16	70.24	1523.20	35.81	9927.68
§5.1		1066.80			1.44			8406.96
§5.2	85.15	116.80	1540.08		0.96	164.32	25.12	2774.59
§5.3				99.12				2031.55
§5.4			70.56					562.03

(c) Results for model 3.6 for Monday logBS in year 1967. This model has no AR(1) component to estimate.

	REML estimation converges in 10 Newton steps						$\hat{\mathbf{y}}$	total
	reduction	$\log \ \mathbf{S}_{\lambda}\ _+$	$\ \hat{\boldsymbol{\beta}}_{\lambda}\ $	$\ \mathbf{V}_{\lambda}\ $	$D_p(\hat{\boldsymbol{\beta}}_{\lambda}, \boldsymbol{\lambda})$	$\log \ \mathbf{H}_{\lambda}\ $		
initial	4046.43	7732.10	23509.83	4356.09	186.91	6043.37	87.88	45962.61
§5.1		2470.95			4.58			40519.13
§5.2	388.07	246.66	8691.42		3.10	662.88	54.49	13902.71
§5.3				504.00				10550.62
§5.4			281.20					2140.40
§5.5	69.30	53.57	61.05	69.41	1.63	108.62	9.73	373.30

## 5.1 A better implementation of the existing computational engine

Inefficiencies in the original computation engine are mainly associated with  $\lfloor \log |S_\lambda|_+ \rfloor$  and  $D_p(\hat{\beta}_\lambda, \lambda)$  steps in REML estimation. They have immediate negative impact on speed as they involve redundant computations.  $\lfloor \log |H_\lambda| \rfloor$  step can be optimized to reduce memory footprint. Its improvement won't be practically observed from execution time, as long as RAM is adequate on the testing platform (which is the case on the Intel Xeon E5-2650 v2 workstation I am working with). See Table 5.2 for the practical effects of these improvements.

### 5.1.1 How $\lfloor \log |S_\lambda|_+ \rfloor$ can be improved

Let us first consider the  $\lfloor \log |S_\lambda|_+ \rfloor$  step. As discussed in §4.3.7, the Frobenius norm  $\lambda_i \|S_i\|_F$  of each individual penalty matrix is computed, and whenever a “bad scaling” between penalty matrices is detected under some numerical tolerance, a similarity transformation is in place. For some reason, implementation in `mgcv` has an interface (see function `gam.reparam`) that only takes the “root”  $S_i^{\frac{1}{2}}$  (see §4.1.5 if you forget what a “root” is) rather than both  $S_i$  and  $S_i^{\frac{1}{2}}$  as input. In consequence,  $S_i$  is re-computed from  $S_i^{\frac{1}{2}}$  at the beginning which could have been avoided. There are also other occasions where using  $S_i^{\frac{1}{2}}$  leads to inefficiency.

1. When similarity transformation  $\dot{U}$  is required,  $S_i$  needs be transformed, too (see §4.3.7). In the general algorithm for  $q > 2$  penalty matrices,  $\dot{U}$  is updated iteratively. `mgcv` chooses to update  $S_i^{\frac{1}{2}}$  along this process, but a more efficient way is to apply a one-off transformation  $\dot{S}_i^{\frac{1}{2}} = \dot{S}_i^{\frac{1}{2}} \dot{U}$  in the end, then compute  $\dot{S}_i = \dot{S}_i^{\frac{1}{2}'} \dot{S}_i^{\frac{1}{2}}$ .
2. When computing second derivatives, `mgcv` computes  $\dot{S}^{-1} \dot{S}_i^{\frac{1}{2}'} \dot{S}_i^{\frac{1}{2}}$  by two matrix multiplications instead of just one multiplication  $\dot{S}^{-1} \dot{S}_i$ . In particular, `mgcv` ignores the triangular structure of  $\dot{S}_i^{\frac{1}{2}}$  (after column permutation), and performs the multiplication as if it is a full dense square matrix. As a result, computations of derivatives are twice as much expensive as what is necessary.

There is yet another re-computation. The  $E_\lambda^*$  (see §4.3.3) required by  $\lfloor \hat{\beta}_\lambda \rfloor$  step can be produced in  $\lfloor \log |S_\lambda|_+ \rfloor$  step. But then `mgcv` re-computes  $\dot{S}$  by  $E_\lambda^{*'} E_\lambda^*$  which is counter-intuitive since it has just obtained  $E_\lambda^*$  from  $\dot{S}$  via a Cholesky factorization!

In summary, the original implementation of  $\lfloor \log |S_\lambda|_+ \rfloor$  step is fairly sub-optimal. Even when a similarity transformation turns out unnecessary, it will still perform some meaningless re-computations. I have modified `gam.reparam` so that it takes both  $S_i$  and  $S_i^{\frac{1}{2}}$  as input. I also choose the more efficient Cholesky method to compute derivatives, instead of the original QR method (see §4.3.2 for details of both methods if you need a revision).

### 5.1.2 How $D_p(\hat{\beta}_\lambda, \lambda)$ can be improved

In `bam`, the five computation steps for REML estimation are actually arranged in a different order (from what I demonstrated in §4.3) as  $\lfloor \log |S_\lambda|_+ \rfloor$ ,  $\lfloor \hat{\beta}_\lambda \rfloor$ ,  $D_p(\hat{\beta}_\lambda, \lambda)$ ,  $\lfloor V_\lambda \rfloor$  and  $\lfloor \log |H_\lambda| \rfloor$ , where  $D_p(\hat{\beta}_\lambda, \lambda)$  step comes before  $\lfloor V_\lambda \rfloor$  step. In this way,  $V_\lambda$  is not available for computing the Jacobian matrix  $J$  (see §4.3.5). Instead, `bam` computes  $J$  solving two triangular systems  $J = -R_\lambda^{*-1} (R_\lambda^{*'}{}^{-1} A)$ , where  $R_\lambda^*$  is the ‘R’ factor from the QR factorization in  $\lfloor \hat{\beta}_\lambda \rfloor$  step (see §4.3.3). Due to the column permutation associated with  $R_\lambda^*$ , rows of  $A$  will be reordered when solving these systems, thus the separation between dense blocks and zero blocks in  $A$  (see Figure 4.10) are destroyed so that the sparsity of  $A$

can not be exploited. As a result, **bam** treats  $\mathbf{A}$  as a dense matrix, making computations much more expensive.

The “fix” is just as what is described in §4.3). Given that  $\mathbf{V}_\lambda$  is required anyway by  $\llbracket \log \mathbf{H}_\lambda \rrbracket$  step, it is more beneficial to schedule its computation ahead to facilitate computations in  $\llbracket D_p(\hat{\beta}_\lambda, \lambda) \rrbracket$  step.

Another inefficiency in  $\llbracket D_p(\hat{\beta}_\lambda, \lambda) \rrbracket$  step relates to the computation of the Hessian matrix  $\frac{\partial^2 D_p}{\partial \rho \partial \rho'}$ . **bam** computes this matrix element-wise using a **for**-loop based on formula (4.7). However, as has been demonstrated in that section, writing the result in a matrix form reveals that many terms will cancel out, and the Hessian matrix can be more efficiently obtained using formula (4.8), which requires no loop at all.

### 5.1.3 How $\llbracket \log \mathbf{H}_\lambda \rrbracket$ can be improved

This is comparatively a minor issue, but can help reduce memory usage during REML estimation. Consider the computation of  $\text{tr}(\mathbf{C}_i \mathbf{C}_j)$  in  $\llbracket \log \mathbf{H}_\lambda \rrbracket$  step. **bam** performs the computation at R-level rather than at C-level. In R, arithmetic operations can not be performed on submatrices directly. (This relates to how “vectorization” is performed in R.)

- For example, the R code `sum(A[1:10, 1:10])` is actually doing `tmp <- A[1:10, 1:10]` and `sum(tmp)`. That is, the submatrix has to be first extracted to a temporary matrix.
- As another example, the R code `A[1:10, 1:10] <- A[1:10, 1:10] + 1` is actually doing `tmp <- A[1:10, 1:10] + 1` and `A[1:10, 1:10] <- tmp`. That is, the replacement submatrix has to be first created as a temporary matrix then be copied into the full matrix.

In Figure 4.11 it has been derived that  $\text{tr}(\mathbf{C}_i \mathbf{C}_j) = \text{tr}(\bar{\mathbf{C}}_i \bar{\mathbf{C}}_j) = \sum_s \sum_t \bar{\mathbf{C}}_i(s, t) \bar{\mathbf{C}}_j'(s, t)$ , where  $\bar{\mathbf{C}}_i$  and  $\bar{\mathbf{C}}_j$  are submatrices of  $\mathbf{C}_i$  and  $\mathbf{C}_j$ . If this computation is done at R-level, there will be too many temporary matrices created (as there is a double loop nest over  $i$  and  $j$ ). In addition, R has to do that matrix transpose explicitly. A better idea is to write compiled code like C code to directly programme this trace computation. For example, if  $\bar{\mathbf{C}}_i$  is in row  $u_1$  to  $u_2$  of  $\mathbf{C}_i$  and  $\bar{\mathbf{C}}_j$  is in row  $v_1$  to  $v_2$  of  $\mathbf{C}_j$ , the trace is just  $\sum_{s=u_1}^{u_2} \sum_{t=v_1}^{v_2} \mathbf{C}_i(s, t) \mathbf{C}_j(t, s)$ , eliminating the need to extract  $\bar{\mathbf{C}}_i$  and  $\bar{\mathbf{C}}_j$ .

## 5.2 Using optimized BLAS for high performance computing

### 5.2.1 An introduction to BLAS

Almost certainly, if you want to do high performance computing, you can’t miss the topic of BLAS.

BLAS (Basic Linear Algebra Subprograms) is a set of standard specifications on a group of elementary numerical operations for common linear algebra computations. It includes the level-1 BLAS (Lawson et al., 1979) for scalar-vector and vector-vector operations, the level-2 BLAS (Dongarra et al., 1988) for matrix-vector operations and the level-3 BLAS (Dongarra et al., 1990) for matrix-matrix operations. These specifications were optimally designed, providing sufficiently flexible functionality that meets practical needs for programming linear algebra computations, while keeping the set of elementary operations as small as possible. For example, the operation **daxpy** from level-1 BLAS does  $\mathbf{y} = \mathbf{y} + \alpha \mathbf{x}$ , that is, multiplying a scalar  $\alpha$  to a vector  $\mathbf{x}$  then adding the resulting vector to  $\mathbf{y}$ . The operation **dgemv** from level-2 BLAS does  $\mathbf{y} = \alpha \text{op}(\mathbf{A})\mathbf{x} + \beta \mathbf{y}$ , where  $\text{op}(\mathbf{A}) = \mathbf{A}'$  or  $\mathbf{A}$  is a matrix,  $\mathbf{x}$  and  $\mathbf{y}$  are vectors,



and  $\alpha$  and  $\beta$  are scalars. The operation `dgemm` from level-3 BLAS does  $\mathbf{C} = \alpha \text{op}(\mathbf{A})\text{op}(\mathbf{B}) + \beta \mathbf{C}$ , where  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{C}$  are matrices, and  $\alpha$  and  $\beta$  are scalars. Together with such specifications, a model implementation written in FORTRAN 77 is provided by Netlib (a repository of software for scientific computing, maintained by AT&T, Bell Laboratories, the University of Tennessee and Oak Ridge National Laboratory). These model subprograms / subroutines / functions further introduce parameters like “stride” and “leading dimension”, to enable operations on vectors whose elements are not necessarily stored contiguously and submatrices from a larger matrix. Such model implementation serves as a stand-alone library, commonly known as the “F77 BLAS” or the “reference BLAS”, and it is portable on all machines.

However, the motivation of BLAS standard, is that BLAS library should be optimized for high performance on any particular machine, by making optimal use of its hardware. The reference BLAS should only serves a sanity check, once an optimized implementation is available. Optimization are often provided by machine’s vendor (but maybe as a proprietary library), like the MKL (Math Kernel Library) (Burylov et al., 2007) from Intel. But a number of open source optimized BLAS libraries are also available and offering competitive performance. These include ATLAS (Automatically Tuned Linear Algebra Software) (Whaley et al., 2000) and OpenBLAS (Zhang et al., 2011) (a project aiming to extend and maintain the successful GotoBLAS (Goto and Geijn, 2008a,b)).

Optimization of BLAS is important and thriving, because BLAS operations, as they were supposed to be, are building blocks for more complicated operations that forms advanced numerical linear algebra libraries. A popular example is the LAPACK (Linear Algebra Package) (Demmel et al., 1987; Anderson et al., 1990, 1999; Anderson and Dongarra, 1990), aiming at solving systems of simultaneous linear equations, least-squares solutions of linear systems of equations, eigenvalue problems, singular value problems, as well as their associated matrix factorizations like LU, Cholesky, QR, SVD, Schur and generalized Schur. With a high-performance BLAS library, LAPACK operations automatically gains efficiency.

All matrix computations introduced in §4.1 can be implemented with existing level-3 BLAS subroutines or LAPACK routines. For demonstration, let  $\mathbf{A}$  and  $\mathbf{B}$  be square matrices,  $\mathbf{S}$  be a positive-definite matrix,  $\mathbf{R}$  be a full-rank upper triangular matrix,  $\mathbf{P}$  be a column permutation matrix and  $\mathbf{Q}$  be the ‘Q’ factor (in factored form storage) from a QR factorization, then

- the matrix cross-product  $\mathbf{A}'\mathbf{A}$  (or  $\mathbf{A}\mathbf{A}'$ ) can be computed using level-3 BLAS subroutine `dsyrk`;
- the matrix multiplication  $\mathbf{AB}$  (or  $\mathbf{A}'\mathbf{B}$ ,  $\mathbf{AB}'$ ,  $\mathbf{A}'\mathbf{B}'$ ) can be computed using level-3 BLAS subroutine `dgemm`;
- triangular matrix multiplication  $\mathbf{RB}$  (or  $\mathbf{R}'\mathbf{B}$ ) can be computed using level-3 BLAS subroutine `dtrmm`;
- solving a triangular system of equations  $\mathbf{R}^{-1}\mathbf{B}$  (or  $\mathbf{R}'^{-1}\mathbf{B}$ ) can be computed using level-3 BLAS subroutine `dtrsm`;
- pivoted QR factorization  $\mathbf{AP} = \mathbf{QR}$ , can be computed using LAPACK routine `dgeqp3`;
- orthonormal transformation  $\mathbf{Q}'\mathbf{B}$  can be computed using LAPACK routine `dormqr`;
- pivoted Cholesky factorization  $\mathbf{P}'\mathbf{SP} = \mathbf{R}'\mathbf{R}$  can be computed using LAPACK routine `dpstrf`;
- positive-definite matrix inverse  $\mathbf{S}^{-1}$  can be computed by LAPACK routine `dpotri`, which consists of a call to `dtrtri` for triangular matrix inverse and a call to `dlaum` for transposed triangular matrix cross-product;
- symmetric eigen decomposition can be computed using LAPACK routine `dsyevr`.

As with many scientific software, R relies on BLAS and LAPACK to do linear algebra computations. For example, the R function “%\*%” for matrix multiplication is mapped to the level-3 BLAS subroutine

`dgemm`. Particularly, BLAS library is a shared library that is dynamically linked to R and automatically loaded at R's startup time. Such dynamic linking can be altered at any time by a user, so it is straightforward to experiment R programs with different BLAS libraries to compare performance. Very often, R is linked to the reference BLAS library, because many users (including me) don't have an existing BLAS library on their machines prior to R installation, thus R will create a reference BLAS library for its use. Therefore, linking R to an optimized BLAS library, like MKL and OpenBLAS, are very likely to yield performance boost to a program, if it is rich in BLAS and LAPACK operations.

### 5.2.2 Benchmarking BLAS

After installing an optimized BLAS library, the first thing to do is benchmarking. In this section, I will compare the performance of the matrix computations listed above in R, using reference BLAS, OpenBLAS and MKL (if available). Hardware information on tested machines is given in Appendix A. Here, performance is reported in GFLOPs (short for GigaFLOPs). 1 GFLOPs =  $10^9$  FLOPs, and FLOPs means "floating-point operations per second". Apart from symmetric eigen decomposition, the exact FLOP count for all matrix computations in the list is explicitly known given dimension of the matrices. The performance can then be derived by dividing FLOP count by measured computation time.

Note that an optimized BLAS library is often multi-threaded to exploit multi-core hardware for parallel computing, but there are several different ways to control the number of threads accessible to the library. To start with, I will pin the number of threads at 1 for strict serial computing.

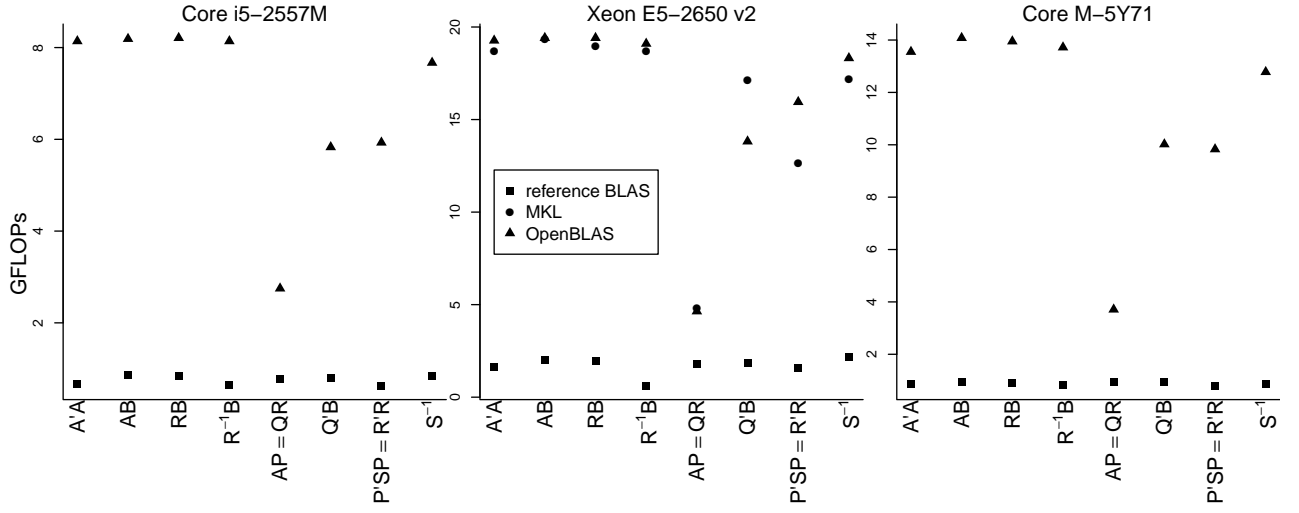
Figure 5.1 illustrates the result of my simple benchmarking (symmetric eigen decomposition is excluded as its exact FLOP count is unknown). From the figure, the following observations can be made.

- On all machines, optimized BLAS is more than 10 times faster than the reference version for level-3 BLAS operations. On the Intel Xeon E5-2650 v2 workstation where both MKL and OpenBLAS are available, these two libraries seem to deliver similar performance.
- Performance of LAPACK routines varies between routines. Positive-definite matrix inverse has competitive performance with level-3 BLAS subroutines. Performance of pivoted Cholesky factorization and orthonormal transformation is lower but still good. But pivoted QR factorization has the poorest performance.
- Performance of LAPACK routines also depends on the version of BLAS. On the Intel Xeon E5-2650 v2, MKL outperforms OpenBLAS for orthonormal transformation, but the other way round for pivoted Cholesky factorization.

In the next two sections, I will explain why optimized BLAS is much faster than the reference BLAS, and why LAPACK routines differ in performance.

### 5.2.3 Block algorithms and data caching

The performance advantage of an optimized BLAS to the reference BLAS comes from the block algorithms it uses for matrix computations to achieve data caching. In scientific computing, using block algorithms is a primary step for high-performance computing as this gives the best cache performance (Lam et al., 1991).



**Figure 5.1:** BLAS and LAPACK's single-threading performance (in GFLOPs)

Our computers have many data storage layers. The closest layer to CPU is a set of CPU registers. CPU can access data on these units with no delay (a more appropriate term is “latency”). Registers are where CPU performs computations. All data must be brought to registers before CPU can process them. However, registers are very small in numbers and size. Even holding a  $4 \times 4$  matrix is impossible. So data must be frequently exchanged between registers and other memory layers during computations: those requested by the current instruction are read in, and those not to be used within the next few instructions will be kicked out. A CPU cache, or L1 (level-1) data cache, is a secondary memory layer right under the registers. It has a bigger size (normally 32KB) and is also fast enough. You can put in a double-precision floating-point vector of 4096 elements, or put in a  $64 \times 64$  matrix, or a few smaller matrices (say three  $36 \times 36$  matrices). Below this CPU cache there may be other cache layers, like L2 cache; then below all cache layers is the main RAM. During computation, data will be frequently transported between all layers. When CPU requests a datum, the farther it is from registers, the longer time it takes to fetch it to registers. Furthermore, different layers have different access speed (or latency). It takes at most 1 CPU cycle to access a register, then normally 4 cycles to access L1 data cache, then maybe 20 cycles to access L2 data cache, and finally, 80 to 200 cycles to access main RAM. In consequence, how fast a program can run depends on how fast data CPU can obtain data it demands. Algorithms with different data access patterns generally give different performance.

For example, to compute a matrix multiplication  $C = AB$  between two  $1200 \times 1200$  square matrices  $A$  and  $B$ , consider the following two arrangements.

1. Compute 1200 matrix-vector multiplications  $C(:, j) = AB(:, j)$ , updating  $C$  column by column;
2. Partition each matrix into 40 small  $30 \times 30$  blocks like

$$A = \begin{pmatrix} A_{1,1} & A_{1,2} & \cdots & A_{1,40} \\ A_{2,1} & A_{2,2} & \cdots & A_{2,40} \\ \vdots & \vdots & \ddots & \vdots \\ A_{40,1} & A_{40,2} & \cdots & A_{40,40} \end{pmatrix}$$

and do block matrix-matrix multiplication  $C_{I,J} = \sum_{K=1}^{40} A_{I,K} B_{K,J}$  to update  $C$  block-by-block.

In the first case, each iteration over  $j$ -loop needs to access the whole matrix  $A$ . A  $1200 \times 1200$  matrix is far too large to be retained in either registers or L1 data cache. After  $C(:, j)$  is updated, most elements of  $A$  would have been evicted from L1 cache, so that  $A$  needs to be loaded from lower memory layers for the update of  $C(:, j+1)$ . The second case appears to be less efficient as it involves more

levels of loop, however, L1 cache is large enough to hold  $\mathbf{A}_{I,K}$ ,  $\mathbf{B}_{K,J}$  and  $\mathbf{C}_{I,J}$ . There will be no data eviction from L1 cache to lower memory layers during computation of  $\mathbf{C}_{I,J}$ . If you make a careful count, you will find that in the first case, each element of  $\mathbf{A}$  is read from memory layers below L1 cache 1200 times, while in the second case, this value is only 40. In the second case, CPU can get data it needs faster, so the matrix-matrix multiplication is faster.

The second algorithm is just an example of a block algorithm. It has an effect of “data caching” (each element of  $\mathbf{A}$  is reused 30 times in L1 cache before being kicked out) and has a higher performance. See Drepper (2007, §3) and Hennessy and Patterson (2011, Chapter 2) for more about CPU cache, and Kowarschik and Weiß (2003) for an overview of block algorithms and data caching in numerical computations.

The above simplified matrix-matrix multiplication example just explains why an optimized BLAS is faster than the reference BLAS. An optimized BLAS implements block algorithms for level-3 BLAS operations. They have a `dgemm`-based structure, by highly optimizing matrix-matrix multiplication and rebuilding other level-3 BLAS operations using the optimized `dgemm`. Such design idea dates back to Kågström et al. (1998), then was superseded by van de Geijn and Quintana-Ortí (2008, Chapter 5) and Gunnels et al. (2001), and was recently reinterpreted by Van Zee and van de Geijn (2015) and Van Zee et al. (2016).

In a block algorithm, there is one or more parameters called *blocking factors*. In the above matrix-matrix multiplication example, 30 is a (and the only) blocking factor. Blocking factors affect the effectiveness of caching. Too large a value destroys caching while too small a value underuses cache capacity.

It is worth pointing out that within BLAS, only level-3 operations can substantially benefit from block algorithms. The purpose of a block algorithm is to achieve data reuse and that reuse takes place in L1 cache. Level-1 BLAS for vector-vector operations has no data reuse at all as all vectors are only visited once. Level-2 BLAS for matrix-vector operations has no data reuse for the matrix operand. The vector operand may be reused but the gains in performance won’t be substantial, because the computational complexity of a matrix-vector operation is determined by the dimension of the matrix. For example, the optimized matrix-vector multiplication from OpenBLAS only has 2.9 GFLOPs performance on an Intel Xeon E5-2650 v2, but the matrix-matrix multiplication attains 19.41 GFLOPs performance. In fact, it is exactly this difference between three levels of operations that makes LAPACK routines differ greatly in performance.

#### 5.2.4 Understanding the performance difference between LAPACK routines

Block algorithms are also used for matrix factorizations, although they are generally more complicated. Such algorithms are constantly developed for LAPACK and *LAPACK Working Notes* has collected hundreds of research work on the theme. LAPACK aims to reorganize basic computational algorithms to use level-3 BLAS operations as much as possible, ending up with block matrix computations in the innermost loop. In this way, the high performance of an optimized BLAS can be carried as much as as possible.

In §4.1 I have explained algorithms for some matrix computations. They can be easily programmed using level-1 and level-2 BLAS operations. However, these basic algorithms are primarily used for educational purposes and are not those behind optimized level-3 BLAS or LAPACK. For some examples of block algorithms, see Quintana-Ortí et al. (1998) for block pivoted QR factorization, Bischof and Loan (1987) for block orthonormal transformation and Lucas (2004) for block pivoted Cholesky factorization. For positive-definite matrix inverse, the basic algorithm for triangular matrix inverse in Figure 4.3a can be extended to a block version straightforward, treating all elements like  $R_{ii}$  as block matrices. Rearranging the algorithm in Figure 4.3b for transposed triangular matrix cross-product

to be block-oriented is less straightforward. I am unable to locate an exact reference for the block implementation in LAPACK, so I encourage all keen readers to read the source code of `dlauum`.

Although a LAPACK routine is intended to be rich in level-3 BLAS operations, some routines still consist of a significant proportion of level-2 BLAS operations (primarily the matrix-vector multiplication). It turns out that for block pivoted QR factorization, there are still exactly 50% of its computations being matrix-vector multiplications. Let  $r_2$  and  $r_3$  respectively be the performance of matrix-vector and matrix-matrix multiplications, then pivoted QR factorization should have an overall performance of  $2/(\frac{1}{r_2} + \frac{1}{r_3})$ . On an Intel Xeon E5-2650 v2 with OpenBLAS, It can be measured that  $r_2 \approx 2.9$  and  $r_3 \approx 19$ , so pivoted QR factorization is expected to have a performance of 5 GFLOPs, which agrees with measured performance in Figure 5.1.

By contrast, pivoted Cholesky factorization is richer in matrix-matrix multiplications. Let  $b$  and  $t$  respectively be the blocking factor and the dimension of the matrix, then matrix-matrix multiplication accounts for  $(1 - \frac{3b}{2t})$  fraction of its computations (Lucas, 2004). LAPACK would choose  $b = 64$  (see `ilaenv` of LAPACK), so for example for an  $5000 \times 5000$  matrix, this proportion is 98.08% which is almost 100%.

### 5.2.5 Using optimized BLAS for GAM computations

Following §4.2 and §4.3, it is clear that GAM computations are rich in the matrix computations I have just benchmarked. Thus, using an optimized BLAS library is very promising to speed up model fitting. From Figure 5.1 it appears that OpenBLAS generally outperforms MKL on an Intel Xeon E5-2650 v2 workstation (see Appendix A for hardware information), so I will link R with OpenBLAS for testing.

Experiment results are presented in Table 5.2, under row “§5.2”. Compared with previous row “§5.1”, the speedup is not uniform for all computation steps, but mostly agrees with what I have explained in previous sections.

- The acceleration of pseudo QR reduction is most dramatic, as its computations is dominated by matrix cross-product for which an optimized BLAS library attains the best performance.
- For REML estimation, the performance improvement on  $[\hat{\beta}_\lambda]$  is limited as it is rich in pivoted QR factorization. Otherwise for  $[\log |S_\lambda|_+]$  and  $[\log |H_\lambda|]$  steps, the gains is substantial.

However, it is surprising to see that  $[V_\lambda]$  has no performance improvement at all. Why?

## 5.3 Computing $V_\lambda$ with BLAS

It turns out that in `bam`, computation of  $V_\lambda$  is programmed with BLAS-free C code, thus changing BLAS has no effects on it. As is listed in §5.2.1, there is an existing LAPACK routine `dpotri` for such computation, so I proceed to use this routine. The experiment result in Table 5.2 proves this very worthwhile.

The C routines in `bam` are inferior to `dpotri` in the following aspects.

1. Triangular matrix inverse  $\tilde{R} = R^{-1}$  is computed by solving triangular systems, as is covered in §4.1.4. But `dpotri` uses the algorithm in Figure 4.3a with a block implementation for caching.

2. Computation of  $\tilde{\mathbf{R}}\tilde{\mathbf{R}}'$  does not exploit the triangular structure of  $\tilde{\mathbf{R}}$ , as is suggested by the algorithm in Figure 4.3b. The resulting matrix multiplication costs  $p^3$  FLOP rather than  $\frac{1}{3}p^3$  FLOP.

## 5.4 Computing $\hat{\beta}_\lambda$ via Cholesky factorization

After previous optimization,  $\hat{\beta}_\lambda$  is now the slowest part hence the bottleneck of GAM computations.

As is described in 4.3.3,  $\hat{\beta}_\lambda$  step is dominated by a pivoted QR factorization of a  $2p \times p$  augmented model matrix. However, it is already observed in §5.2 that such factorization can not attain high performance with an optimized BLAS library. So a good idea is to find an alternative algorithm than using pivoted QR factorization.

In fact, this is possible. From (1.8) we see that  $\hat{\beta}_\lambda$  solves a system of linear equations  $\mathbf{H}_\lambda \beta = \mathbf{X}'\mathbf{W}\mathbf{y}$  with a positive definite or positive-semidefinite coefficient matrix  $\mathbf{H}_\lambda = \mathbf{X}'\mathbf{W}\mathbf{X} + \mathbf{S}_\lambda$ . Now consider a pivoted Cholesky factorization with pre-conditioning (see §4.1.2 if you need a revision)  $\mathbf{H}_\lambda = (\mathbf{R}_\lambda^*)' \mathbf{R}_\lambda^*$  where  $\mathbf{R}_\lambda^* = \mathbf{R}_\lambda \mathbf{R}_\lambda' \mathbf{J}_\lambda$ , then  $\hat{\beta}_\lambda$  can be efficiently obtained by solving two triangular systems (see §4.1.3 if you need a revision)  $\tilde{\mathbf{z}} = \mathbf{R}_\lambda^{*-1} \mathbf{z}$  and  $\hat{\beta}_\lambda = \mathbf{R}_\lambda^{-1} \tilde{\mathbf{z}}$ . Such computation only involves  $(\frac{1}{3}p^3 + 2p^2)$  FLOP, much lower than the  $\frac{10}{3}p^3$  FLOP using the QR factorization method in §4.3.3. In addition, that pivoted Cholesky factorization outperforms pivoted QR factorization in an optimized BLAS implies that the practical gains from using this Cholesky method can be more than what the FLOP count shows. The results in Table 5.2 shows that it is 20 times faster for model 3.5 and 30 times faster for 3.6!

Computations of  $D_p(\hat{\beta}_\lambda, \lambda)$  and  $\log |\mathbf{H}_\lambda|$  are also convenient in this method. There are  $D_p(\hat{\beta}_\lambda, \lambda) = \|\mathbf{W}^{\frac{1}{2}}\mathbf{y}\|^2 - \|\tilde{\mathbf{z}}\|^2$  and  $\log |\mathbf{H}_\lambda| = 2 \sum_{i=1}^p (\log |\mathbf{R}_\lambda(i, i)| + \log \mathbf{J}_\lambda(i, i))$ .

In case  $\mathbf{H}_\lambda$  is rank-deficient with rank  $r < p$ , computations proceed as follows. Compute  $\mathbf{b} = \mathbf{P}_\lambda' \mathbf{J}_\lambda^{-1} \mathbf{z}$ . Initialize  $\tilde{\mathbf{z}}$  as a length- $p$  vector of zeros and set  $\tilde{\mathbf{z}}(1:r) = \mathbf{R}_\lambda(1:r, 1:r)'^{-1} \mathbf{b}(1:r)$ . Reset  $\mathbf{b}(1:r) = \mathbf{R}_\lambda(1:r, 1:r)^{-1} \tilde{\mathbf{z}}(1:r)$  and  $\mathbf{b}((r+1):p) = \mathbf{0}$ , then compute  $\hat{\beta}_\lambda = \mathbf{J}_\lambda^{-1} \mathbf{P}_\lambda' \mathbf{b}$ ,  $D_p(\hat{\beta}_\lambda, \lambda) = \|\mathbf{W}^{\frac{1}{2}}\mathbf{y}\|^2 - \|\tilde{\mathbf{z}}(1:r)\|^2$  and  $\log |\mathbf{H}_\lambda| = 2 \sum_{i=1}^r \log |\mathbf{R}_\lambda(i, i)| + 2 \sum_{i=1}^p \log \mathbf{J}_\lambda(i, i)$ .

Note that the matrix cross-product  $\mathbf{X}'\mathbf{W}\mathbf{X}$  and  $\mathbf{z} = \mathbf{X}'\mathbf{W}\mathbf{y}$  required by this method are readily accumulated in pseudo QR reduction (see Figure 4.6 if you need to remaster the algorithm); In addition, the final Cholesky factorization in the original pseudo QR reduction will no longer be necessary.

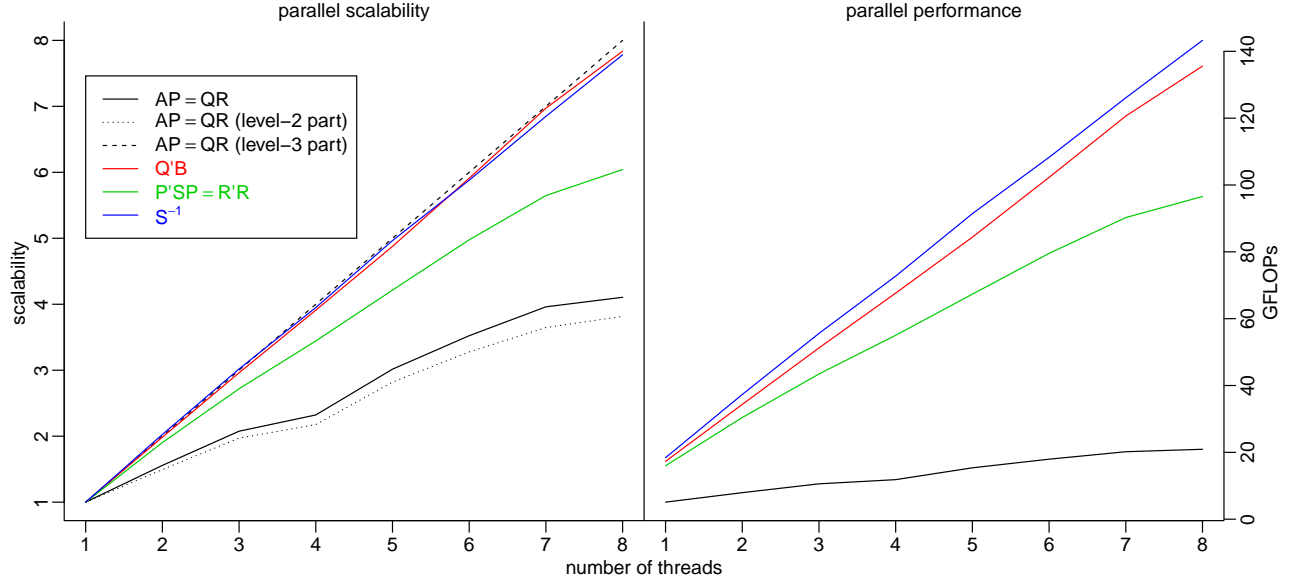
## 5.5 Using parallel computing for GAM computations

Previously I have been using OpenBLAS for serial computing, now I will investigate further speedup via parallel computing.

### 5.5.1 Parallel performance of BLAS and LAPACK

Let us start by testing the parallel performance of BLAS and LAPACK. As it turns out, the parallel performance for level-3 BLAS operations has almost perfect parallel scaling (see §4.2.3 if you need a

revision on parallel speedup and parallel scaling): using  $m$  CPU cores makes the computations almost  $m$  times as fast. So more interest is on comparing parallel performance between LAPACK routines.



**Figure 5.2:** Parallel scaling (left panel) and parallel performance (right panel, unit: GFLOPs) for LAPACK operations: pivoted QR factorization (black), orthonormal transformation (red) pivoted Cholesky factorization (green) and positive-definite matrix inverse (blue). Particularly, in the left panel, the scaling for level-2 (black dotted) and level-3 (black dashed) BLAS computations inside pivoted QR are also displayed. Experiment is conducted on an Intel Xeon E5-2650 v2 workstation with OpenBLAS (see Appendix A for hardware information). All matrices used for the experiment are  $N_m \times N_m$  square matrix, but the dimension grows with the number of threads  $m$  with relationship  $N_m = \sqrt[3]{m}N_1$ . Since all operations involve  $O(N_m^3)$  FLOP, this makes sure the loading for each thread is constant at  $O(N_1^3)$  no matter how  $m$  grows. For this benchmarking, I have chosen  $N_1 = 10000$ . The parallel scaling factor measured this way is known as the *weak parallel scaling factor*, since the problem size grows with  $m$ . Another type of parallel scaling factor is *strong parallel scaling factor*, where the problem size is held fix. That is, using  $N_m = N_1$ . Strong scaling is an ideal property, but more difficult to achieve than weak scaling. This is not difficult to understand. If problem size is fixed, then the loading for each thread will diminish as the number of cores increases, while at the same time the overhead for parallel scheduling will increase.

Figure 5.2 illustrates measured parallel performance of the four LAPACK routines listed in §5.2.1. The left panel sketches parallel scaling factor against the number of CPU cores  $m$ , and the right panel sketches the measured performance (in GFLOPs) against  $m$ . There are the following observations.

- Positive-definite matrix inverse (blue line) and orthonormal transformation (red line) have near perfect parallel performance;
- The parallel scaling of pivoted Cholesky factorization (green line) is inferior, but still remarkable. There is a 6 times speedup when 8 cores are used;
- The parallel scaling of pivoted QR factorization (black solid line) is poor. Using 8 cores only achieves 4 times speedup, and the scaling curve begins plateaus.

So why does pivoted QR factorization scale so poorly? The reason is still with the large proportion of level-2 BLAS operations it involves. As is mentioned earlier, parallel scaling for level-3 BLAS operations is impressively perfect. However, the scaling for level-2 BLAS operations is not even as half as good. See the left panel of Figure 5.2 again, and look for the black dashed line and the black dotted line. The former is the parallel scaling (good) for the level-3 BLAS computations inside pivoted QR factorization, and the latter is the parallel scaling (poor) for the level-2 BLAS computations. Since level-2 BLAS operations account for 50% of computations of pivoted QR factorization, it limits the parallel scaling of pivoted QR factorization. Other LAPACK operations have demonstrated good parallel scaling, because they involve a smaller proportion of level-2 BLAS operations.

So why does level-2 BLAS scale poorly as opposed to level-3 BLAS? Consider an example of matrix-vector multiplication and matrix-matrix multiplication. For simplicity, all matrices are square of size  $N$ . For large  $N$ , the former performs  $2N^2$  FLOP on  $(N^2 + \frac{N}{b} + N) \approx N^2$  data read and  $N$  data write, while the latter performs  $2N^3$  FLOP on  $(2N^2 + \frac{N^3}{b}) \approx \frac{N^3}{b}$  data read and  $\frac{N^3}{b}$  data write. Here,  $b$  is the blocking factor that MKL or OpenBLAS chooses for caching (see §5.2.3 if you need a revision on blocking factor). The exact derivation of the amount of data motion requires knowledge of the blocking strategy used by optimized BLAS library, but is not our focus here. It is sufficient to observe that the number of FLOP and memory operations is 2 : 1 for the former, but  $b$  : 1 for the latter. On an Intel Xeon E5-2650 v2 workstation with OpenBLAS (see Appendix A for hardware information), MKL and OpenBLAS chooses  $b \approx 256$  (of course, the exact value is machine dependent), so we can conclude that matrix-vector multiplication is *memory-bound* but matrix-matrix multiplication is *computation-bound* or *CPU bound*. Memory-bound operations inflict heavier front-side bus (FSB) pressure, and the FSB bandwidth can easily become a limiting factor for their parallel scalability, which is particularly true in shared-memory parallel computing via multi-threading. With  $N = 10000$ , matrix-vector multiplication takes about 0.069 seconds (hence giving 2.9 GFLOP performance as we claimed earlier), implying that data are transported in FSB at a rate of 11.6 GB/s, which is already 77.7% of a memory channel's capacity. In multi-threading, each thread would demand a data flow at this supplying rate. As soon as the gross demand exceeds the maximum FSB bandwidth of all four channels, which is 59.712 GB/s, the FSB is “jam” and CPU cores have to wait longer for data's arrival. Another implication of this, is that the scalability will be upper bounded by  $59.712/11.6 \approx 5.15$ , however we increase the number of threads, justifying why the black dotted curve in left panel of Figure 5.2 begins to plateau. By contrast, data are only flowing in FSB at 2.53 GB/s for matrix-matrix multiplication. Even if all 16 cores are used, the gross demand is 40.48 GB/s which the FSB has no difficulty to meet. I do make a test with 16 threads, and it turns out that the scaling factor is as high as 15 for matrix-matrix multiplication, but is still held at 4 for matrix-vector multiplication.

A message of this, is that it is a good idea to avoid using pivoted QR factorization for developing algorithms if we wish the algorithm to scale well. Fortunately, the pivoted QR factorization in  $[\hat{\beta}_\lambda]$  step of REML estimation has been replaced by pivoted Cholesky factorization in the last section. In addition, I have been using pseudo QR reduction rather than QR reduction, so no pivoted QR factorization is involved in model matrix reduction step either. These give a very positive message that GAM computations should scale reasonably well in practice. In the next section, I will experiment GAM computations with parallel computing.

### 5.5.2 Experimenting parallel computing for GAM computations

**bam** already has a “fork-join” implementation for parallel pseudo QR reduction, so for this part, I will supply it with single-threaded OpenBLAS. For REML estimation, I will enable multi-threading for OpenBLAS, so that REML estimation will be “automatically” parallelized.

It is generally a good idea to experiment parallel computing on a problem with big size. I will consider the following **test model**:

$$\begin{aligned} \log \text{BS}_{\text{id}} = & f_0(\mathbf{E}_{\text{id}}; 5) + f_1(\mathbf{w}_{\text{y}}; 20) + f_2(\{\mathbf{e}_{\text{i}}, \mathbf{n}_{\text{i}}\}; 200) + f_3(\mathbf{w}_{\text{y}}, \{\mathbf{e}_{\text{i}}, \mathbf{n}_{\text{i}}\}; 20, 200) + \\ & f_4(\mathbf{T}_{\text{id}}^0; 15) + f_5(\mathbf{T}_{\text{id}}^0, \{\mathbf{e}_{\text{i}}, \mathbf{n}_{\text{i}}\}; 15, 200) + \\ & f_6(\mathbf{T}_{\text{id}}^*; 15) + f_7(\mathbf{T}_{\text{id}}^*, \{\mathbf{e}_{\text{i}}, \mathbf{n}_{\text{i}}\}; 15, 200) + \\ & f_8(\mathbf{T}_{\text{id}}^0, \mathbf{T}_{\text{id}}^*; 15, 15) + f_9(\mathbf{i}; 1276) + f_{10}(\mathbf{w}_{\text{y}}; 52) + f_{11}(\mathbf{h}_{\text{i}}; 10) + \\ & f_{12}(\mathbf{d}_{\text{w}}; 7) + f_{13}(\mathbf{d}_{\text{w}}, \mathbf{w}_{\text{y}}; 7, 20) + f_{14}(\mathbf{d}_{\text{w}}, \{\mathbf{e}_{\text{i}}, \mathbf{n}_{\text{i}}\}; 7, 200) + e_{\text{id}}, \end{aligned}$$

This is a very simple extension to model 3.6, as a joint model for logBS from all days of week. The extension procedure is similar to what was previously done in the *Manchester 11* case study. Day of week  $\mathbf{d}_{\text{w}}$  is introduced as a factor variable (levels 1 to 7, for Monday to Sunday), and  $f_1$  and  $f_2$  are allowed to vary for each factor level. To be precise, new components to model 3.6 are



**Table 5.3:** Time (in seconds) for fitting the **test model** assuming i.i.d. model errors (examination of its residuals will be done in the next Chapter). The model has 12460 coefficients for 387253 data. The model is fitted on an Intel Xeon E5-2650 v2 workstation (see Appendix A for hardware information), with R software linked to OpenBLAS. Up to 8 CPU cores are used for parallel computing. For pseudo QR reduction, the “fork-join” parallelism is implemented with single-threaded OpenBLAS; for REML estimation, multi-threading is enabled for OpenBLAS. Also see Figure 5.3 for illustration of parallel scaling.

$m$	reduction	REML estimation converges after 22 iterations					$\hat{\mathbf{y}}$	total
		$\log \ \mathbf{S}_{\lambda}\ _+$	$\ \hat{\beta}_{\lambda}\ $	$\ \mathbf{V}_{\lambda}\ $	$D_p(\hat{\beta}_{\lambda}, \lambda)$	$\log \ \mathbf{H}_{\lambda}\ $		
1	3675.00	507.8	850.41	1476.02	7.47	1665.56	714.21	8896.47
2	1914.06	286	471.32	746.37	5.09	882.35	371.39	4676.58
3	1318.00	211.05	344.09	513.57	4.35	619.99	249.97	3261.02
4	1011.64	174.06	280.11	397.29	4.02	489.75	192.84	2549.71
5	801.05	152.41	241.73	327.56	3.83	413.56	157.13	2097.27
6	690.91	135.72	216.82	281.78	3.76	358.59	131.41	1818.99
7	599.02	125.46	198.78	248.64	3.72	322.92	116.41	1614.95
8	534.16	118.1	186.03	224.55	3.75	294.47	102.85	1463.91

- $f_{12}(\mathbf{d}_w; 7)$ , an i.i.d. random effect modelling the mean logBS at each day of week;
- $f_{13}(\mathbf{d}_w, \mathbf{w}_y; 7, 20)$ , a tensor product spline, whose first margin is an i.i.d. random effect of  $\mathbf{d}_w$  and the second margin is a cubic cyclic spline of  $\mathbf{w}_y$ . This component models the deviation of logBS seasonality at each day of week from the mean logBS seasonality  $f_1$ ;
- $f_{14}(\mathbf{d}_w, \{\mathbf{e}_i, \mathbf{n}_i\}; 7, 200)$ , a tensor product spline, whose first margin is an i.i.d. random effect of  $\mathbf{d}_w$  and the second margin is a rank 200 thin-plate spline of spatial coordinates. This component models the deviation of logBS over space at each day of week from the mean logBS over space  $f_2$ ;
- $e_{id}$ , model error, assumed to be i.i.d. for model fitting.

Whether this model is reasonable or not is not of interest here. It is chosen for testing just to demonstrate that the new GAM fitting methods developed in this Chapter makes it feasible to estimate a daily logBS model. This test model involves 12460 coefficients for 387253 data.

Table 5.3 presents the fitting time for this model. Using 8 CPU core, model fitting completes in about 25 minutes. Figure 5.3 further illustrates parallel scaling for reduction, REML estimation and the overall GAM fitting. See figure caption for more details. In short, the parallel scaling is very good. (Model fitting time for model 3.6 using 8 CPU cores is also reported in Table 5.2.) In conclusion, the computational hurdle at the start of this Chapter is successfully overcome.

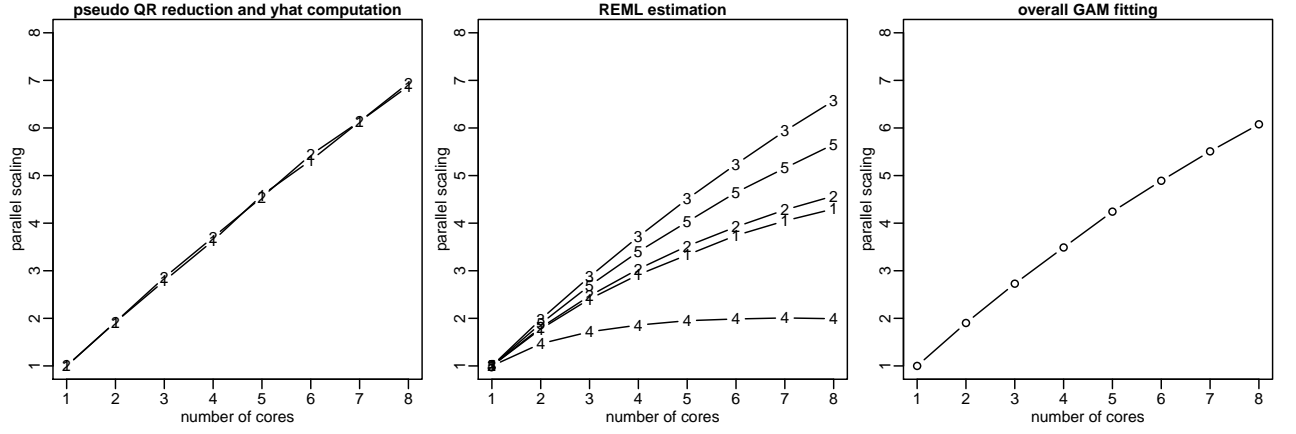
## 5.6 A comparison of computational capability with INLA

GAM estimation in `mgcv` is a type of frequentist shrinkage estimation or empirical Bayes estimation (see §1.1.5). Hierarchical Bayesian estimation treating smoothing parameters (or variance components as they are often known in literature of mixed models) as hyper parameters may be an alternative. The R package `INLA` using integrated nested Laplace approximation provides a practically efficient way for model estimation via this route.

Representing GAMs for logBS in `INLA` is not straightforward, as these models are rich in interaction. As an experiment, I began with an extremely simple GAM with only a single interaction term:

$$\log \text{BS}_{id} = f_1(\mathbf{w}_y) + f_2(\mathbf{d}_w) + f_3(\mathbf{d}_w, \mathbf{w}_y) + f_4(\{\mathbf{e}_i, \mathbf{n}_i\}) + \epsilon_{id}.$$

The following explains how each function is constructed in `mgcv` and `INLA`.



**Figure 5.3:** Illustration of parallel scaling in GAM fitting using methods developed in this Chapter. The **test model** assuming i.i.d. errors is fitted for experiment, and see Table 5.3 for the raw timing statistics used to compute parallel scaling factors in this figure. The panel on the left shows the parallel scaling for “fork-join” pseudo QR reduction (labelled with 1) and  $\hat{\mathbf{y}}$  computation (labelled with 2). They have near equal parallel scaling and are both very impressive. The middle panel shows the parallel scaling for REML estimation, where lines labelled with 1 to 5 are respectively associated with the five computation steps in REML estimation, namely  $\log[\mathbf{S}_{\lambda+}]$ ,  $\hat{\beta}_{\lambda}$ ,  $\mathbf{V}_{\lambda}$ ,  $D_p(\hat{\beta}_{\lambda}, \lambda)$  and  $\log[\mathbf{H}_{\lambda}]$ . The scaling property differs between these steps.  $\log[\mathbf{H}_{\lambda}]$  (“5”) and  $\mathbf{V}_{\lambda}$  (“3”) have very good scaling because these steps primarily involves matrix multiplications.  $\log[\mathbf{S}_{\lambda+}]$  (“1”) and  $\hat{\beta}_{\lambda}$  (“2”) scale slightly inferior, because they need to solve triangular system of linear equations apart from pivoted Cholesky factorization.  $D_p(\hat{\beta}_{\lambda}, \lambda)$  (“4”) scales poorly but this does not matter in practice, since this step takes almost negligible computation time (see for example, Table 5.3). The right panel shows the parallel scaling for overall GAM fitting, which is very good.

- In `mgcv`,  $f_1$  is constructed as a cubic cyclic spline; in `INLA`, it is constructed as a first order random walk.
- $f_2$  is constructed as an i.i.d. random effect in both `mgcv` and `INLA`.
- $f_3$  is constructed as a tensor product spline in `mgcv` between a cubic cyclic spline margin and a random effect margin. In `INLA`, it is constructed as a Kronecker product model (see (Lindgren and Rue, 2015, §3.2)), where there is a first random walk of  $\mathbf{w}_y$  for each level of  $\mathbf{d}_w$ .
- in `mgcv`,  $f_4$  is constructed as a thin-plate spline; in `INLA`, it is constructed as a Gaussian random field with Matérn covariance, represented as a solution to a stochastic partial differential equation (SPDE) (Lindgren et al., 2011).

Fitting this model to all 387253 daily logBS in year 1967 only took 2 minutes in `mgcv` with less than 5 GB RAM, but still did not finish after 3 hours with `INLA`. In particular, `INLA` consumed all 128 GB RAM available on a E5-2650 v2 workstation and began to use increasingly more disk storage for swapping. I terminated the computation as the operating system was brought to a halt. There was then no need to try more complicated models with `INLA`. Fitting daily models with `INLA` is just infeasible.

## Chapter 6

# Preliminary modelling of logBS (a revisit)

Being able to fit a model as fast as possible makes many statistical questions answerable. For example, it is now possible to run cross-validation for model 3.6 and model 3.7 to see whether three-way interactions should be retained in a daily logBS model. In this Chapter, I will revisit daily logBS models from 1967 and annual mean logBS models.

### 6.1 A revisit to daily logBS model in 1967

#### 6.1.1 Justifying three-way interactions with cross-validation

The cross-validation is performed as follows. Randomly partition the Monday logBS dataset (54386 data) into a training subset with 80% of the data and a test subset with the other 20%. Fit model 3.6 (without three-way interactions) and model 3.7 (with three-way interactions) to the training subset, and compute mean prediction squared error (MPSE) of both models on the test subset. Repeat this for 100 times, and compute the average MPSE as well as its 95%-confidence interval for both models. Table 6.1 presents the results. To conclude, these three-way interactions should be retained in the model.

#### 6.1.2 A joint model for all days of week?

Model development so far for daily logBS in 1967 has been restricted to building separate models for each day of week. It is time to ask whether it is possible to build a joint model for all days of week. But a key question is: is this worthwhile? Extending model 3.7 to a joint model by allowing many splines to vary with  $\mathbf{d}_w$  can result in too many parameters, especially if those three-way interactions

**Table 6.1:** Cross-validation test for three-way interactions. “mean” gives the average MPSE over 100 simulations; “lwr” and “upr” respectively give lower and upper bound of the 95% confidence interval for the mean. There is strong evidence that three-way interactions should be retained in the model.

	mean	lwr	upr
model 3.6	0.2289	0.2280	0.2298
model 3.7	0.2167	0.2158	0.2176

are extended to four-way interactions.

Let us first attempt the following test model:

### joint model

$$\begin{aligned} \log \text{BS}_{\text{id}} = & f_0(\mathbf{E}_i; 5) + f_1(\mathbf{w}_y; 10) + f_2(\{\mathbf{e}_i, \mathbf{n}_i\}; 150) + f_3(\mathbf{w}_y, \{\mathbf{e}_i, \mathbf{n}_i\}; 20, 150) + \\ & f_4(\mathbf{T}_{\text{id}}^0; 15) + f_5(\mathbf{T}_{\text{id}}^0, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 100) + \\ & f_6(\mathbf{T}_{\text{id}}^*; 15) + f_7(\mathbf{T}_{\text{id}}^*, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 100) + \\ & f_8(\mathbf{T}_{\text{id}}^0, \mathbf{T}_{\text{id}}^*; 15, 15) + f_9(i; 1275) + f_{10}(\mathbf{w}_y; 52) + f_{11}(\mathbf{h}_i; 5) + \\ & f_{12}(\mathbf{w}_y, \mathbf{T}_{\text{id}}^0, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 10, 30) + f_{13}(\mathbf{w}_y, \mathbf{T}_{\text{id}}^*, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 10, 30) + \\ & f_{14}(\mathbf{d}_w; 7) + f_{15}(\mathbf{d}_w, \mathbf{w}_y; 7, 20) + f_{16}(\mathbf{d}_w, \{\mathbf{e}_i, \mathbf{n}_i\}; 7, 150) + \\ & f_{17}(\mathbf{d}_w, \mathbf{T}_{\text{id}}^0; 7, 10) + f_{18}(\mathbf{d}_w, \mathbf{T}_{\text{id}}^0; 7, 10) + f_{19}(\mathbf{d}_w, \mathbf{w}_y, \{\mathbf{e}_i, \mathbf{n}_i\}; 7, 10, 50) + e_{\text{id}}. \end{aligned}$$

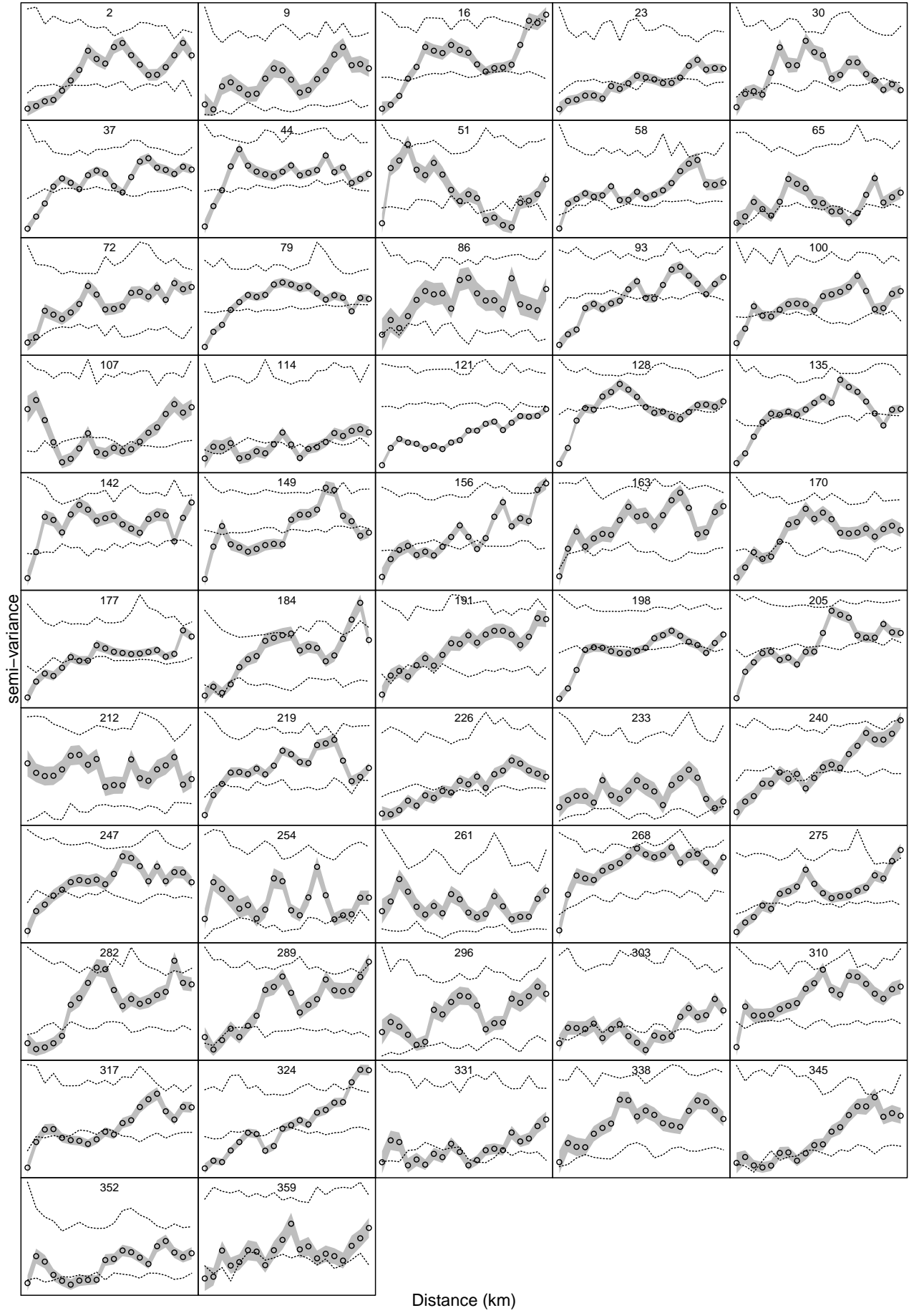
New components to model 3.7 are

- $f_{14}(\mathbf{d}_w; 7)$ , an i.i.d. random effect modelling the mean logBS at each day of week;
- $f_{15}(\mathbf{d}_w, \mathbf{w}_y; 7, 20)$ , a tensor product spline, whose first margin is an i.i.d. random effect of  $\mathbf{d}_w$  and the second margin is a cubic cyclic spline of  $\mathbf{w}_y$ . This component models the deviation of logBS seasonality at each day of week from the mean logBS seasonality  $f_1$ ;
- $f_{16}(\mathbf{d}_w, \{\mathbf{e}_i, \mathbf{n}_i\}; 7, 200)$ , a tensor product spline, whose first margin is an i.i.d. random effect of  $\mathbf{d}_w$  and the second margin is a rank 200 thin-plate spline of spatial coordinates. This component models the deviation of logBS over space at each day of week from the mean logBS over space  $f_2$ ;
- $f_{17}(\mathbf{d}_w, \mathbf{T}_{\text{id}}^0; 7, 10)$ , a tensor product spline, whose first margin is an i.i.d. random effect of  $\mathbf{d}_w$  and the second margin is a cubic regression spline of  $\mathbf{T}_{\text{id}}^0$ . This component models how the average effect of daily minimum temperature  $f_4$  may vary within a week;
- $f_{18}(\mathbf{d}_w, \mathbf{T}_{\text{id}}^*; 7, 10)$ , a similar component to  $f_{17}$  associated with diurnal temperature difference;
- $f_{19}(\mathbf{d}_w, \mathbf{w}_y, \{\mathbf{e}_i, \mathbf{n}_i\}; 7, 10, 50)$ , a three-margin tensor product spline modelling how the spatially varying seasonality  $f_3$  can further vary within a week;
- $e_{\text{id}}$ , model error. This is definitely not i.i.d.. An appropriate assumption may be proposed after examining residuals of this model fitted under i.i.d. error assumption.

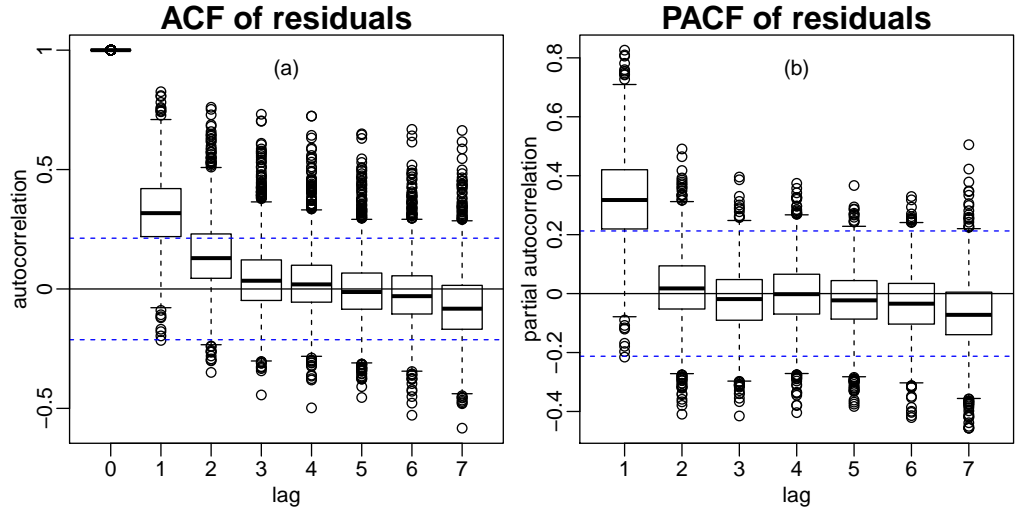
Note that the  $k$  value for the spatial margin in  $f_{19}$  is deliberately kept low to avoid ending up with too many parameters. Still, this model has 14394 regression coefficients.

A key concern is how good this joint model is in modelling space-time relation. It is straightforward to produce variograms of its spatial residuals on all 365 days, but it is probably sufficient to just produce such variograms at all 52 Mondays, so that a comparison can be made with Figure 3.46 from the fitted Monday-only model. This is done in 6.1, and it turns out that the joint model is very inadequate in spatial modelling. Just looking at the first row, residuals of the joint model have unmodelled spatial autocorrelation on day 2, 16, 23, 30, while residuals of the Monday-only model on these days are all fine.

The model may be improved if  $f_{12}$  and  $f_{13}$  are extended to four-way interactions. However, this will raise the number of coefficients to 39450. Even if there is no computational difficulty, why should we fit this model at all? If all components of model 3.7 should vary with day of week, then why not just build seven separate models?



**Figure 6.1:** Empirical variograms of spatial residuals from fitted joint model on each of the 52 Mondays in year 1967. Compared with Figure 3.46 from fitted Monday-only model, the joint model is very inadequate in spatial model.



**Figure 6.2:** ACF (panel (a)) and PACF (panel (b)) of residuals from all stations in fitted joint model (assuming i.i.d. errors). The results are presented in boxplots, and the blue dashed lines are “95% confidence interval” computed based on the median sample size over all stations. The cutoff of PACF at lag 1 shows that on average, residuals from each station are AR(1) correlated.

Building separate models for each day of week does offer some convenience in practice. On one hand, temporal autocorrelation won’t step in model development; on the other hand, the number of data used for model fitting is always 1 / 7 of the volume of the complete dataset hence model fitting is always faster.

However, as part of methodology development, it is useful to demonstrate how temporal autocorrelation may be dealt with in a model for daily data. In the next section, I will explain the construction of block AR(1) errors where the AR(1) coefficient can be estimated by the golden-section search method previously covered. And in the next next section I will give some practical tricks that help speed up the golden-section search when fitting GAMs for large datasets.

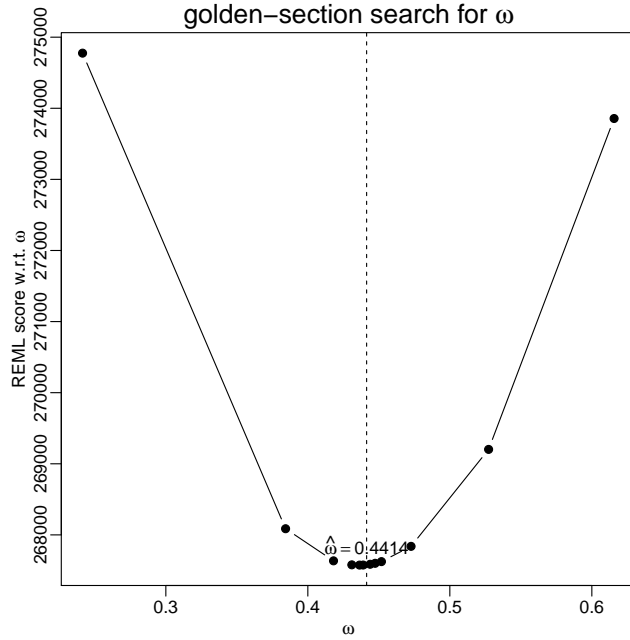
### 6.1.3 Block AR(1) errors with autocorrelation coefficient $\omega$

Fitting the joint model assuming i.i.d. error leaves unmodelled temporal autocorrelation in residuals, as panel (a) of Figure 6.2 shows. Panel (b) also shows that the average pattern of such autocorrelation can be modelled by an AR(1) process. In reality, the autocorrelation coefficient varies greatly across stations (and I have previously demonstrated this in the histogram from Figure 3.24), but estimating all these individual coefficients is not practically feasible. For example, the golden-section search method introduced in §2.2.3 can only deal with a single unknown coefficient. So here the simplest strategy is taken, by only modelling the average autocorrelation coefficient. In other words, errors  $e_{i,d}$  from station  $i$  are an AR(1) process

$$e_{i,d} = \omega \cdot e_{i,d-1} + \epsilon_{i,d-1},$$

where the autocorrelation coefficient  $\omega$  is common for all stations. This produces a block diagonal weight matrix for the resulting additive model, where each diagonal block is a symmetric tri-diagonal matrix like (2.1). For example, if there are two AR(1) processes, each of length 4, then the weight matrix is

$$W = \begin{pmatrix} A & B & & & & & \\ B & C & B & & & & \\ & B & C & B & & & \\ & & B & A & 0 & & \\ & & & 0 & A & B & \\ & & & & B & C & B \\ & & & & & B & C & B \\ & & & & & & B & A \end{pmatrix}, \quad (6.1)$$



**Figure 6.3:** Illustration of golden-section search in estimation of  $\omega$  for the joint model. It converges in 14 steps, giving a point estimate  $\hat{\omega} = 0.4414$ .

where  $A = 1/(1 - \omega^2)$ ,  $B = -\omega A$  and  $C = (1 + \omega^2)A$ . There are two ways to store this matrix: storing each diagonal separately or storing its non-zero elements by row. Either way, elements from the lower or upper sub-diagonal need not be stored due to symmetry (this explains why the lower sub-diagonal in the above matrix is coloured grey). For standardization (see §2.2.2 if you need a revision), the upper triangular Cholesky factor of this weight matrix is

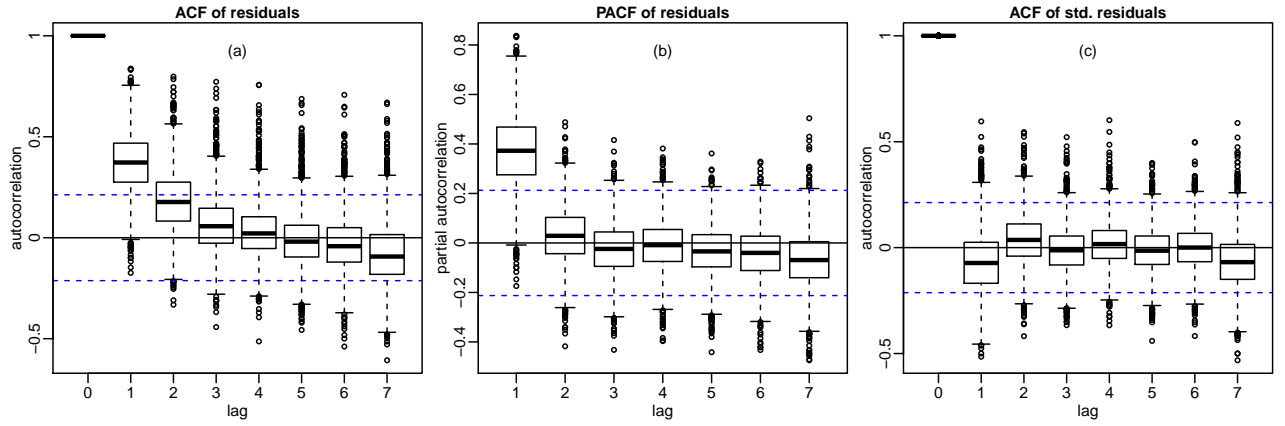
$$\mathbf{W}^{\frac{1}{2}} = \begin{pmatrix} D & E & & & & & & \\ & D & E & & & & & \\ & & D & E & & & & \\ & & & 1 & 0 & & & \\ & & & \hline & & & D & & & & \\ & & & & E & & & \\ & & & & D & E & & \\ & & & & & D & E & \\ & & & & & & 1 & \end{pmatrix},$$

where  $D = \sqrt{A}$  and  $E = -\omega\sqrt{A}$ . To estimate  $\omega$  via golden-section search, the log-determinant of the weight matrix is required for evaluation of REML score. The block structure makes this computation very easy, as the determinant (or log-determinant) of the full weight matrix is just the product (or sum) of the determinant (or log-determinant) for each block. In §2.2.2 I have shown that the log-determinant of a length- $l_i$  AR(1) process is  $-(l_i - 1)\log(1 - \omega^2)$ , thus it is straightforward to see that  $\log|\mathbf{W}| = -(n - \gamma)\log(1 - \omega^2)$ , where  $\gamma$  is the number of AR(1) blocks and  $n = \sum_i l_i$  is the total number of data from all blocks.

Figure 6.3 illustrates the golden-section search for the joint model. It converges in 14 steps, giving a point estimate  $\hat{\omega} = 0.4414$ . Figure 6.4 sketches the ACF and PACF of model residuals, as well as the ACF of standardized residuals. It can be seen that on average, residual autocorrelation are eliminated.

#### 6.1.4 Speeding up golden-section search

Golden-section search makes GAM fitting substantially more costly. With 14 steps for convergence, it is essentially fitting the joint model for 14 times, so the overall fitting time becomes as long as 7 hours. In practice, the fitting process can be slightly speeded up by specifying a narrower search interval.



**Figure 6.4:** ACF (panel (a)) and PACF (panel (b)) of residuals from all stations in fitted joint model (assuming AR(1) errors). Panel (c) further gives the ACF of standardized residuals. The results are presented in boxplots, and the blue dashed lines are “95% confidence interval” computed based on the median sample size over all stations. The AR(1) assumption for mean errors over all stations is adequate, and on average, standardized residuals at each station are uncorrelated.

Previously a very general interval (0.01, 0.99) has been used, now we may restrict it to, for example (0.3, 0.5). As a result, the search converges in 10 steps. This reasonable guess is made by spotting Figure 6.2, where the mean lag-1 autocorrelation is slightly above 0.3. This value is always lower than the one to be estimated via golden-section search, because some proportion of autocorrelation has been fitted by splines under i.i.d. error assumption. So it is safe to use 0.3 as the lower bound of the search interval. 0.5 is chosen as the upper bound because this value is readily much greater than the 75% quantile of lag-1 autocorrelation.

Another way to speed up convergence is to use a slightly higher numerical tolerance (see Figure 2.12) for convergence test. Previously  $\epsilon = 0.001$  has been used, but this is probably more than necessary. In particular, since the objective function is flat near its minimum (the first derivative is close to 0), search steps are increasingly smaller near minimum and the decrease in function value is insignificant. There isn’t much point spending hours on this tiny improvement. Increasing the tolerance to 0.005 requires only 6 steps for convergence, given a point estimate at 0.4413 which is basically no different from 0.4414.

A even better way for faster convergence is to use an adaptive tolerance for REML estimation (see  $\epsilon$  in Figure 4.9). Previously for each trial  $\omega$  from the outer golden-section loop, 17 Newton-Raphson steps are taken to achieve convergence subject to  $\epsilon = 10^{-6}$ . Think about this carefully, there is little sense to iterate Newton-Raphson to convergence with such a high precision before the golden-section has a sign for convergence. Of course, while this is a reasonable idea, practically programming this needs some care to ensure stability of the algorithm. At the moment I haven’t implemented this idea.

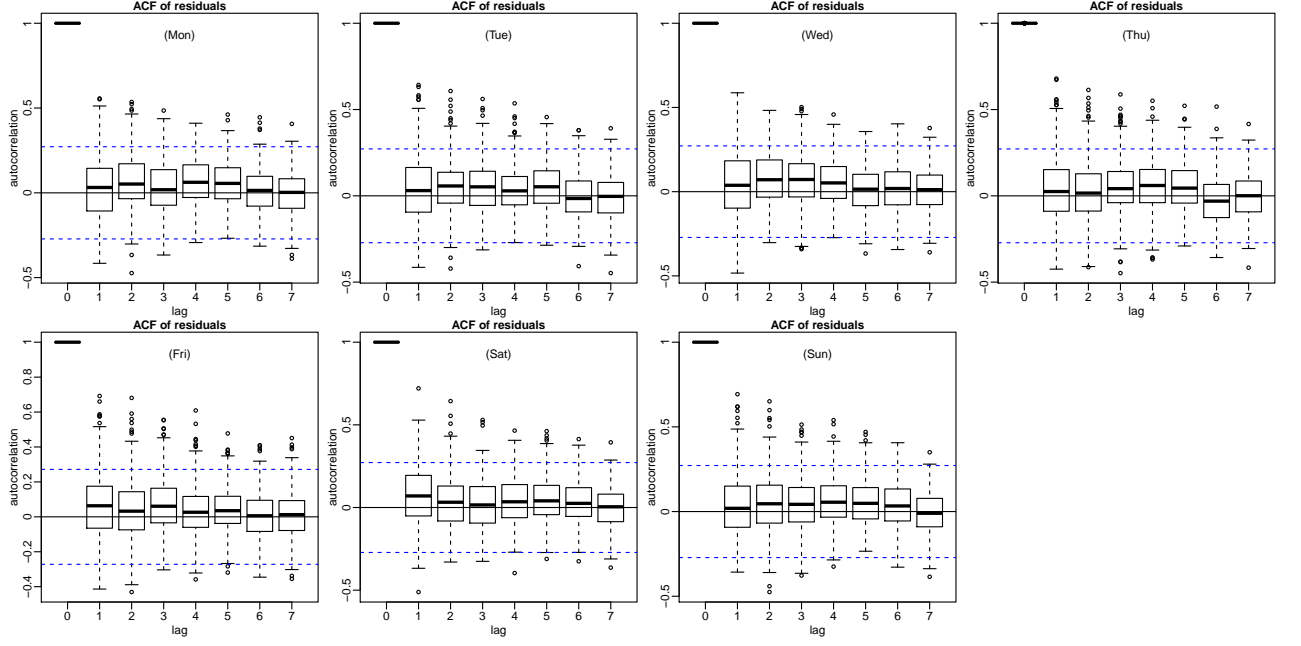
### 6.1.5 Summary and visualization of seven separate models

The seven fitted models have already been summarized in Table 3.10. I will now show some diagnostic plots for these models.

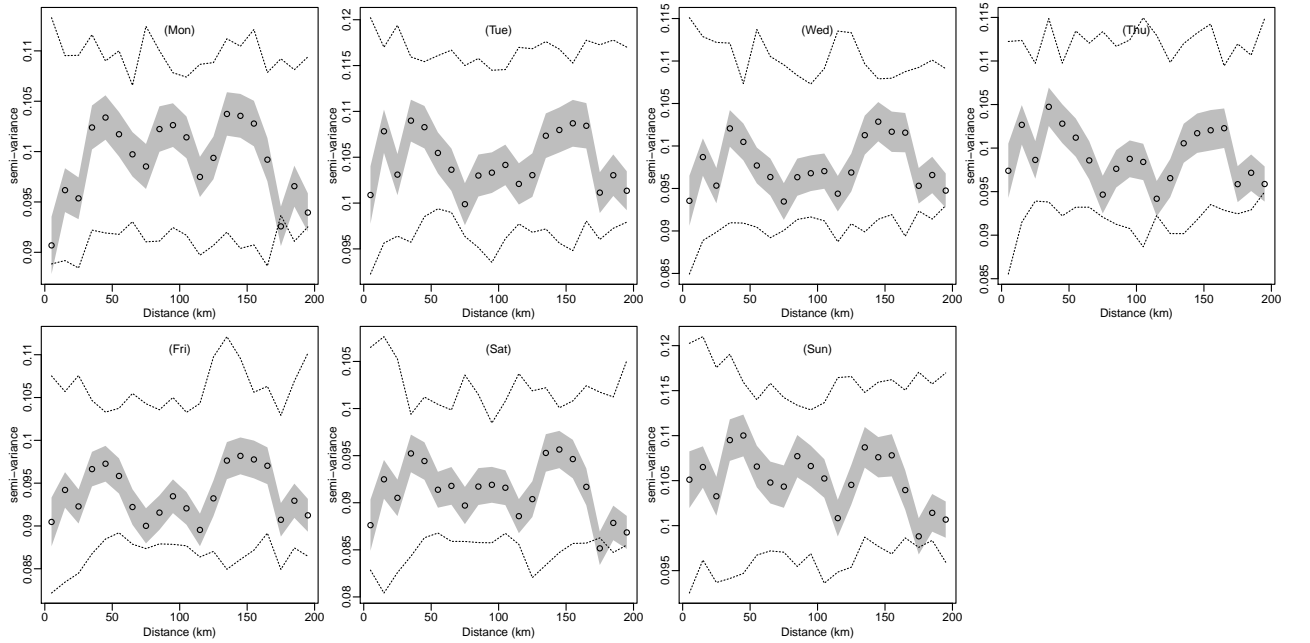
Figures 6.5, 6.6 and 6.7 respectively check i.i.d. assumption for model errors, station-specific random effect and week of year random effect. No violation is seen. Figure 6.8 sketches residuals against fitted values on each day of week. I am not surprised that there seems to be some mean-variance dependence at low fitted values and high fitted values; this is just a reflection of the difficulty to model extreme values in daily logBS. Overall, there is no violation of constant error variance assumption.

Without a joint model it is less straightforward to compare components between different models.

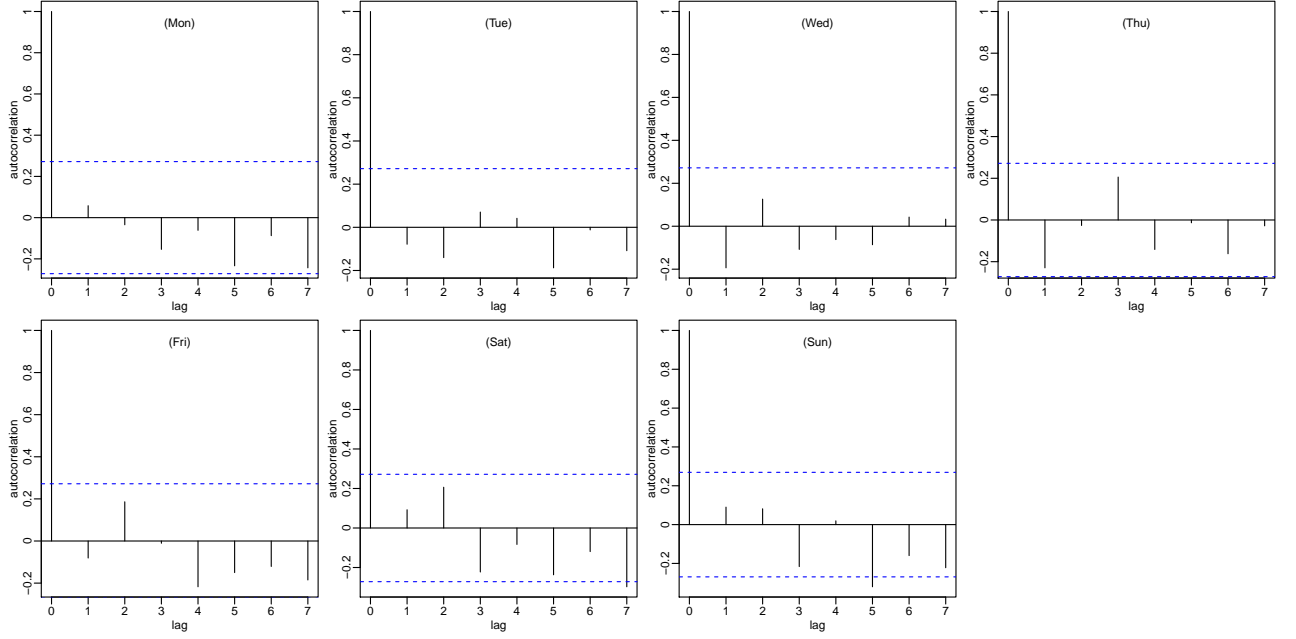




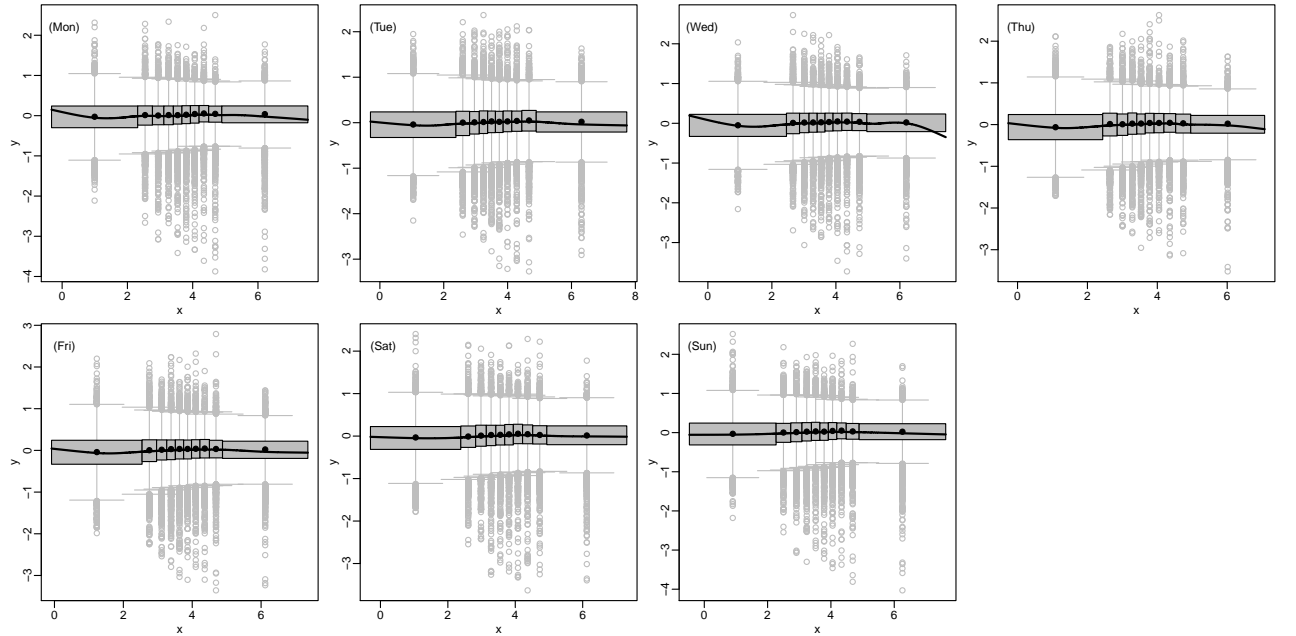
**Figure 6.5:** ACF of residuals of the fitted model 3.6 for each day of week. Results are presented in boxplots, and the blue dashed lines are “95% confidence interval” computed based on the median sample size over all stations. It is evident that on average, there is no temporal autocorrelation in residuals, so the i.i.d. model error assumption is justified.



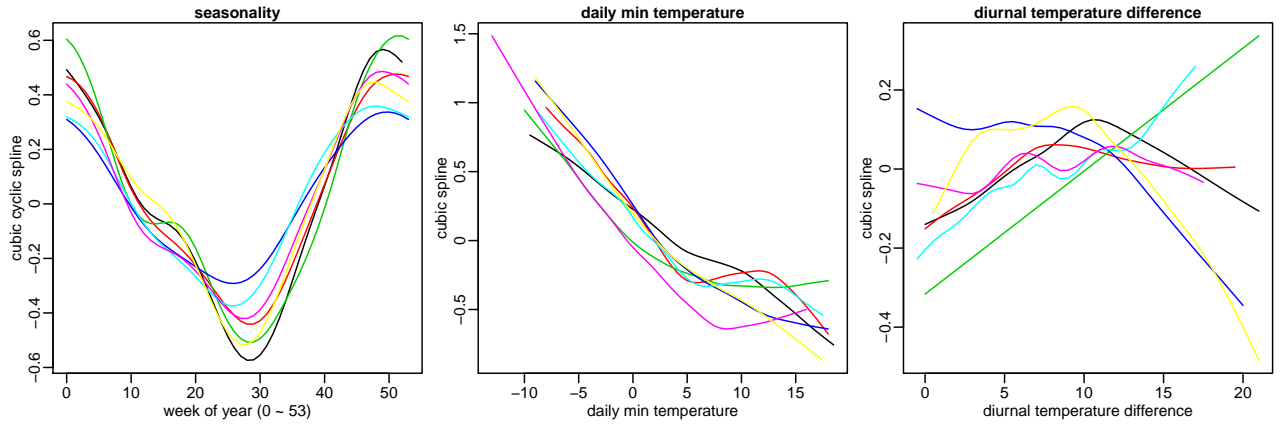
**Figure 6.6:** Variogram of station-specific i.i.d. random effect  $\hat{f}_9$  in the fitted model 3.6 for each day of week. No spatial autocorrelation is spotted, thus the i.i.d. assumption is validated.



**Figure 6.7:** ACF of week of year i.i.d. random effect  $\hat{f}_{10}$  in the fitted model 3.6 for each day of week. No temporal autocorrelation is spotted, thus the i.i.d. assumption is not violated.



**Figure 6.8:** Residuals against fitted values for each day of week. I am not surprised that there seems to be some mean-variance dependence at low fitted values and high fitted values; this is just a reflection of the difficulty to model extreme values in daily logBS. Overall, there is no violation of constant error variance assumption.



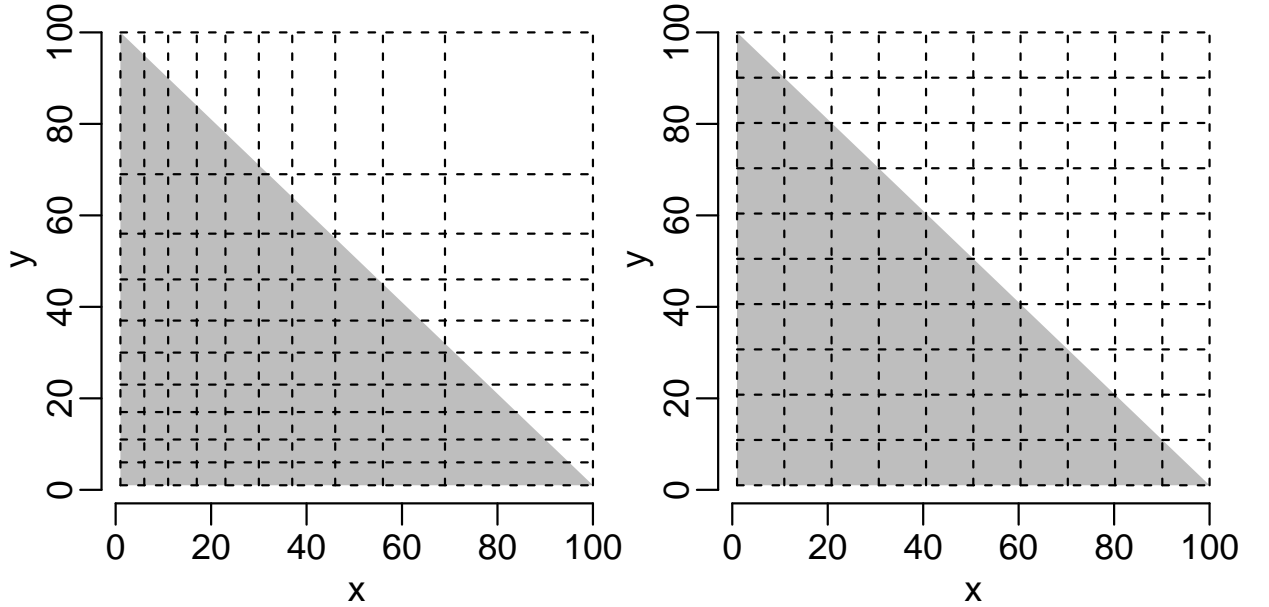
**Figure 6.9:** Estimated seasonality and main temperature effects from different models (black: Monday; red: Tuesday; green: Wednesday; blue: Thursday; cyan: Friday; magenta: Saturday; yellow: Sunday). Note that these curves are not comparable in values, as different models have different intercepts and those curves are subject to different centring constraints. Only their shapes can be compared.

It is also difficult to assess weekly effects. So Figure 6.9 comparing estimated seasonality and main temperature effects from different models is very informal. The only useful message is probably that these effects have a different shape on each day of week. The effects of diurnal temperature difference is very difficult to interpret. I think this once again signals that more daily covariates are needed for building models, otherwise relationship between logBS and those unavailable daily variables have to be described by sophisticated relationship between logBS and temperature.

### 6.1.6 Summary

In this section I revisited the spatial-temporal modelling of daily logBS in year 1967, completing the previously unfinished discussion regarding three-way interactions between space, time and temperature. The cross-validation result supports the inclusion of these components in a daily model. However, the presence of such three-way interactions makes it difficult to extend separate models per day of week to a joint model. A relatively modest joint model without extending those three-way interactions to four-way interactions was attempted, but it turned out inadequate in spatial modelling. It was then concluded that building an adequate joint model would result in a model with too high complexity, so I finally stayed with the strategy of building separate models.

Building separate models for each day of week eliminates the needs for modelling temporal autocorrelation in daily data, but I still discussed how such autocorrelation can be dealt with, by assuming model errors at each station to be AR(1) with the same correlation coefficient  $\omega$ . This assumption makes estimation of  $\omega$  straightforward via golden-section search. Perhaps the assumption of a common  $\omega$  is over restricted. In principle it is fairly easy to assume a different  $\omega_i$  for each station: the resulting weight matrix is still symmetric tri-diagonal, with  $A_i$ ,  $B_i$  and  $C_i$  in place of  $A$ ,  $B$  and  $C$ . However, this ends up with too many parameters in the weight matrix and how to estimate so many  $\omega_i$ 's? A probably more moderate strategy is to assume that  $\omega_i$  varies smoothly over space. In fact, thinking toward this direction suggests that  $\mathbf{W}^{-1}\phi$  be constructed as a spatial-temporal covariance matrix. The link can be established as follows. A closer look at (6.1) reveals that the weight matrix is just the Kronecker product (“ $\otimes$ ”) between an identity matrix  $\mathbf{I}$  and the inverse correlation matrix of an AR(1) process  $\mathbf{W}_t$  (as in (2.1)). If  $\mathbf{I}$  is thought of as the inverse correlation of an i.i.d. spatial process, then it is certainly natural to generalize it to an inverse spatial correlation matrix  $\mathbf{W}_s$  and construct  $\mathbf{W} = \mathbf{W}_s \otimes \mathbf{W}_t$ . The complexity of  $\mathbf{W}$  is well under control, because  $\mathbf{W}_t$  just has a single parameter  $\omega$  and  $\mathbf{W}_s$  often just has a single parameter as well. However, to make practical computation (as well as storage) feasible,  $\mathbf{W}_s$  has to be very sparse. Implementation of this in GAMs could be an interesting topic to pursue in future.



**Figure 6.10:** A bivariate tensor product spline example to illustrate different knots placement strategies. Suppose we want to estimate a tensor product spline  $s(x, y)$  defined on a  $[1, 100] \times [1, 100]$  square, but data are only sampled at grid points  $(i, j)$ ,  $i, j = 1, 2, \dots, 100$  in the lower triangular domain (the gray shaded area). If both margins are modelled by a cubic regression spline with 11 knots, the left panel places knots at 0%, 10%, ..., 100% quantiles on each margin, while the right panel places knots evenly on each margin. In both panels, intersections of horizontal and vertical dashed lines give the resulting knots of the tensor product spline.

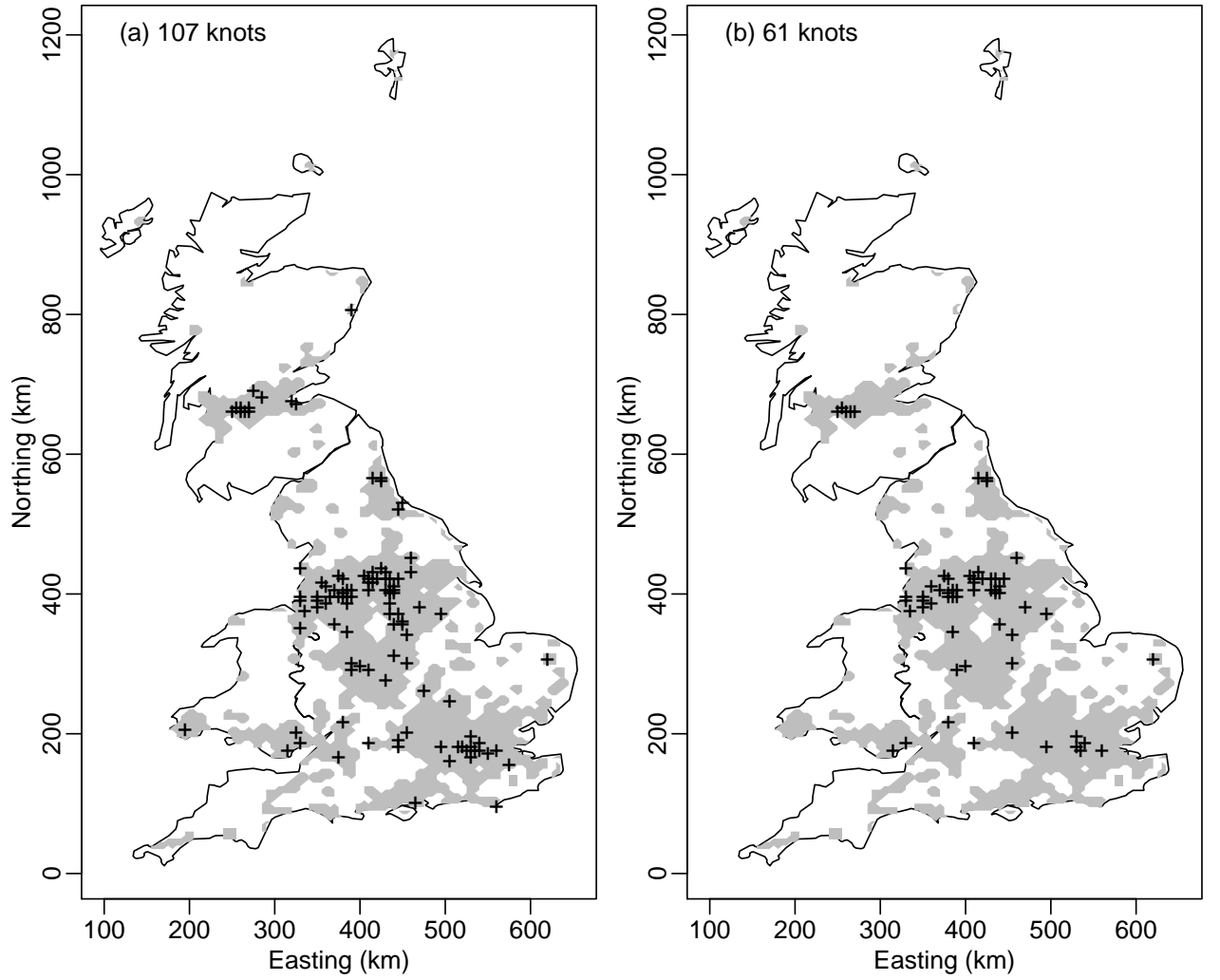
## 6.2 A revisit to annual mean logBS model

Previously in §3.4.6 I have realized the extrapolation problem in making spatial-prediction as the Network became sparser and sparser. Now I am going to make an attempt to address this problem. My attempt turned out unsuccessful, but it is still interesting to present it here.

From Figure 3.34 it is obvious that the tensor product spline  $f_3(\mathbf{y}, \{\mathbf{e}_i, \mathbf{n}_i\})$  modelling the space-time interaction is the cause of extrapolation. A conjecture is that the thin-plate spline margin is overflexible in later years when there were insufficient spatial information. So I began to think about a way to restrict flexibility of this margin.

The issue may be related to knots placement. To motivate this, consider a bivariate tensor product spline example illustrated in Figure 6.10. Knots placement in the left panel is more sensible, as knots are more densely distributed in the sampling region. Fewer knots in the unsampled region means less risk of extrapolation over there. Analogy can be made for  $f_3(\mathbf{y}, \{\mathbf{e}_i, \mathbf{n}_i\})$ . Since I am to model space-time interaction with a tensor product spline, knots should probably be placed in space-time coordinates where information of such interaction is rich. In other words, knots of the thin-plate spline should be placed at spatial locations where long-term Black Smoke sampling is available. However, the automatic knots placement for a thin-plate spline in `mgcv` disregards this, ending up using all unique 2626 site locations. Had all stations had an adequately long monitoring history, such auto knots placement would not be problematic, but stations of the Black Smoke Network generally had a very short life (see Figure 3.25). Therefore, such auto knots placement may result in unreasonable spatial-temporal extrapolation.

To obtain a sensible set of knots, consider those  $5\text{km} \times 5\text{km}$  pixels of the spatial prediction domain (see Figure 3.1). Among all 3949 pixels, 107 have observations for at least 30 years and 61 have observations for at least 35 years. Panels (a) and (b) of Figure 6.11 respectively pinpoint these pixels, which are respectively suitable as knots of a rank-107 and a rank-61 thin-plate spline.



**Figure 6.11:** Attempts to mitigate spatial-temporal extrapolation by placing knots of the thin-plate spline at spatial locations where long-term Black Smoke sampling is available. Panels (a) and (b) respectively select 125, 75 and 30 pixels from the spatial domain (see Figure 3.1) as knots, where at least 30 and 35 years' observations are available.

**Table 6.2:** Cross-validation for different knots placement strategies for thin-plate spline margin of  $f_3$  in model 3.5. Column “mean” gives average MPSE from 100 simulations; “lwr” and “upr” respectively give lower and upper bound of 95%-confidence interval of “mean”.

	$m = 1$			$m = 2$			$m = 3$		
	mean	lwr	upr	mean	lwr	upr	mean	lwr	upr
auto	0.6282	0.6238	0.6326	0.6255	0.6216	0.6294	0.5856	0.5824	0.5888
panel (a)	0.6266	0.6220	0.6311	0.6264	0.6228	0.6300	0.5966	0.5933	0.5998
panel (b)	0.6151	0.6106	0.6195	0.6137	0.6102	0.6171	0.5824	0.5791	0.5856
	$m = 4$			$m = 5$					
	mean	lwr	upr	mean	lwr	upr			
auto	0.5745	0.5706	0.5784	0.5668	0.5633	0.5704			
panel (a)	0.5851	0.5808	0.5894	0.5742	0.5706	0.5779			
panel (b)	0.5758	0.5717	0.5798	0.5671	0.5633	0.5709			

To justify such manual knots placement, I decided to do cross-validation which is designed as follows. Randomly select 1000 stations and delete their last  $m$  years' observations. Use the deleted data as test dataset and the rest of the data as training dataset. Fit model 3.5 to the training dataset, with three different strategies for knots placement: auto, panel (a), panel (b), then compute the mean prediction squared error on test dataset. The simulation is repeated for 100 times. In addition, I will try  $m = 1, 2, 3, 4, 5$ . Table 6.2 presents the results of cross-validation. It is hard to draw a firm

conclusion from those results. Strategy “panel (a)” almost always gives higher prediction error than “auto”, except for  $m = 1$  case. Then strategy “panel (b)” gives lower prediction error than “auto” for  $m = 1, 2, 3$ , but the prediction error is then higher “auto” for  $m = 4, 5$ . So there is no solid evidence that manual knots placement yields a better model. Unfortunately I have to admit that I can not resolve this issue in this thesis.

# Chapter 7

## A model for all daily logBS?

Plenty of logBS models have been built by now, but only a small proportion of the full daily logBS dataset has been used. What if I am to fit a model using the complete dataset? Is this computationally feasible?

### 7.1 A test model

The following is a test model for daily logBS on Monday:

**Model 7.1:**

$$\begin{aligned} \log\text{BS}_{\text{id}} = & f_0(\mathbf{E}_i; 5) + f_1(\mathbf{h}_i; 5) + f_2(\{\mathbf{e}_i, \mathbf{n}_i\}; 150) + f_3(\mathbf{i}; 2874) + \\ & f_4(\mathbf{w}; 40) + f_5(\mathbf{w}; 2310) + f_6(\mathbf{E}_i, \mathbf{w}; 5, 40) + f_7(\mathbf{w}, \{\mathbf{e}_i, \mathbf{n}_i\}; 40, 150) + \\ & f_8(\mathbf{w}_y; 10) + f_9(\mathbf{w}_y, \{\mathbf{e}_i, \mathbf{n}_i\}; 20, 150) + f_{10}(\mathbf{w}_y, \mathbf{w}; 20, 40) + \\ & f_{11}(\mathbf{T}_{\text{id}}^0; 15) + f_{12}(\mathbf{T}_{\text{id}}^*; 15) + f_{13}(\mathbf{T}_{\text{id}}^0, \mathbf{T}_{\text{id}}^*; 15, 15) + \\ & f_{14}(\mathbf{T}_{\text{id}}^0, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 100) + f_{15}(\mathbf{T}_{\text{id}}^*, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 100) + \\ & f_{16}(\mathbf{w}_y, \mathbf{T}_{\text{id}}^0, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 10, 30) + f_{17}(\mathbf{w}_y, \mathbf{T}_{\text{id}}^*, \{\mathbf{e}_i, \mathbf{n}_i\}; 10, 10, 30) + \\ & f_{18}(\mathbf{T}_{\text{id}}^0, \mathbf{w}; 10, 40) + f_{19}(\mathbf{T}_{\text{id}}^*, \mathbf{w}; 10, 40) + f_{20}(\mathbf{r}_{\text{im}}; 10) + \epsilon_{\text{iw}}. \end{aligned}$$

- $f_0$  is a random effect modelling the mean of logBS from stations of environmental type  $\mathbf{E}_i$  (5 types in total);
- $f_1$  is a natural cubic spline of  $\mathbf{h}_i$  with 5 knots;
- $f_2$  is a rank-150 thin-plate spline modelling the general spatial distribution of logBS across stations (invariant with time);
- $f_3$  is an i.i.d. random effect for station  $\mathbf{i}$  (2874 stations in total);
- $f_4$  is a natural cubic spline modelling the general long term trend (common to all stations);
- $f_5$  is an i.i.d. random effect for  $\mathbf{w}$  (2310 weeks in total);
- $f_6$  is a tensor product spline, whose first margin is a random effect for  $\mathbf{E}_i$  and the second margin is cubic spline of  $\mathbf{w}$ . This component models how  $f_4$  varies for each  $\mathbf{E}_i$ ;
- $f_7$  is a tensor product spline, whose first margin is a natural cubic spline of  $\mathbf{w}$  and the second margin is a thin-plate spline of  $\{\mathbf{e}_i, \mathbf{n}_i\}$ . This component models how  $f_4$  varies over space;

- $f_8$  is a cubic cyclic spline of  $\mathbf{w}_y$  with 20 knots, modelling the mean seasonality of logBS at all stations;
- $f_9$  is a tensor product spline whose first margin is a cubic cyclic spline of  $\mathbf{w}_y$  and the second margin is a thin-plate spline of  $\{\mathbf{e}_i, \mathbf{n}_i\}$ . This component models how  $f_8$  varies over space;
- $f_{10}$  is a tensor product spline whose first margin is a cubic cyclic spline of  $\mathbf{w}_y$  and the second margin is a natural cubic spline of  $\mathbf{w}$ . This component models how  $f_8$  varies in a long term;
- $f_{11}$  is a natural cubic spline with 10 knots for daily minimum temperature, modelling the mean effect of this variable at all stations;
- $f_{12}$  is a similar component to  $f_{11}$ , associated with diurnal temperature difference;
- $f_{13}$  is a tensor product spline modelling the interaction between two temperature variables;
- $f_{14}$  is a tensor product spline whose first margin is a natural cubic spline of  $T_{id}^0$  and the second margin is a thin-plate spline of  $\{\mathbf{e}_i, \mathbf{n}_i\}$ . This component models how  $f_{11}$  varies over space;
- $f_{15}$  is a similar component to  $f_{14}$ , associated with diurnal temperature difference;
- $f_{16}$  and  $f_{17}$  are three-way interactions (three-margin tensor product splines) between space, time and temperature. These components model how  $f_{14}$  and  $f_{15}$  vary within a year;
- $f_{18}$  and  $f_{19}$  are two-margin tensor product splines, modelling how  $f_{11}$  and  $f_{12}$  vary in a long term;
- $f_{20}$  is natural cubic spline of monthly rainfall, modelling the mean effect of this variable at all stations in all time;
- $\epsilon_{i\mathbf{w}}$  is an i.i.d. model error.

This model is proposed based on previous experience on logBS modelling.

- Components  $f_0$  to  $f_7$  are motivated by model 3.5 for annual mean logBS, except that time variable year  $y$  is replaced by week  $\mathbf{w}$ , motivated by model 3.1 in *Manchester 11* case study.
- Components  $f_8$  to  $f_{13}$  are also motivated by *Manchester 11* case study, where it was shown that time-varying seasonality could be modelled by an interaction between  $\mathbf{w}_y$  and  $\mathbf{w}$  and that temperature variables are useful for predicting logBS;
- Components  $f_{14}$  to  $f_{17}$  are motivated by model 3.7 for Monday logBS in 1967, where it was shown that temperature effects would vary over space and time;
- Components  $f_{18}$  and  $f_{19}$  are added to model how relationship between logBS and temperature might change over decades (In §3.2.6 I mentioned that such relationship can not be estimated from *Manchester 11* case study due to small sample size);
- $f_{20}$  is added for monthly rainfall. This variable was not included in any previous models, because the relationship between daily logBS and monthly rainfall appeared very weak (see for example Figures 3.13 and 3.40). However, adding a low-rank cubic spline for rainfall does not make model estimation any more difficult; if this variable indeed has no effect, this spline may be “penalized out”.

Note that I can go on adding more components. For example, by allowing  $f_{16}$  and  $f_{17}$  to vary with  $\mathbf{w}$ , or by allowing  $f_{18}$  and  $f_{19}$  to vary with space. However, this model already has 22571 regression coefficients. There are in total 1342159 Monday logBS observations, so on average, there is one coefficient per 59.5 data.

This test model may be useful for two purposes.



1. The difficulty of modelling daily logBS can be assessed. The test model is readily flexible enough, and if it turns out inadequate, it is very likely that logBS can not be properly modelled using GAMs and there is no further point to go on building a joint model for all days of week.
2. It serves a very good example for investigating GAM computation methods. More than half of the model components are tensor product splines and there are even three-margin tensor product splines. Fitting such model is not a trivial task. Questions may be raised on whether the GAM computation methods developed in Chapter 5 should be further optimized to accommodate such complicated models.

## 7.2 Fitting the test model

How long does it take to fit the test model using methods in Chapter 5? It is interesting to first make some prediction about this.

- Pseudo QR reduction is rich in matrix cross-product  $\mathbf{X}'\mathbf{W}\mathbf{X}$ , which an optimized BLAS is best at. On an Intel Xeon E5-2650 v2 workstation (see Appendix A for hardware information), OpenBLAS attains a performance of 250 GFLOPs using all 16 CPU cores for matrix-matrix multiplication (revisit §5.2.2 if you have forgotten what “GFLOPs” is). The test model has  $n = 1342159$  data and  $p = 22571$  coefficients. As is summarized in §4.2.4, pseudo QR reduction involves  $np^2$  FLOP, that is 683763 GFLOP for the test model. Thus I predict that the reduction should takes about  $683763 / 250 = 2735$  seconds, which is about 45 minutes.
- As is summarized in §4.3.8, the FLOP count for each computation step is different. In particular, many matrix computations of difference performance are involved, so computation time in REML estimation is much more difficult to estimate. A simple yet reliable prediction is to run the REML estimation for just a single Newton-Raphson step, because the total time needed by REML estimation is just this time multiplied by the number of iterations. To start REML estimation, results of pseudo QR reduction, namely  $\mathbf{X}'\mathbf{W}\mathbf{X}$  and  $\mathbf{X}'\mathbf{W}\mathbf{y}$  are required. However, there is no need to really do the reduction. For performance measurement, it is sufficient to use a  $p \times p$  random matrix for  $\mathbf{X}'\mathbf{W}\mathbf{X}$  and a length- $p$  random vector for  $\mathbf{X}'\mathbf{W}\mathbf{y}$ . It turns out that for the test model, a single Newton step takes 146 seconds on an Intel Xeon E5-2650 v2 workstation with 16-threaded OpenBLAS. Usually REML estimation takes 10 to 25 steps for convergence, so REML estimation for the test model is unlikely to take more than an hour.

Estimation above assured me that fitting the test model is feasible. So I proceeded to actually fit this model. However, practical computation took more time than what I had expected. It turned out that

- REML estimation took 21 steps for convergence, spending 45 minutes, which agrees with my estimate;
- Pseudo QR reduction took 1 hour and 15 minutes, much longer than my estimate;
- Computation of  $\hat{\mathbf{y}}$  took about 25 minutes.

So, I forgot to take into account of the computation of fitted values, residuals, etc. But, why is  $\hat{\mathbf{y}}$  computation so expensive? Is it not a matrix-vector multiplication which in principle only involves  $2np$  FLOP? On an Intel Xeon E5-2650 v2 workstation matrix-vector multiplication has a performance of 2.9 GFLOPs, so even in serial computing this should just take about 21 seconds. What is contributing to that 25 minutes?

**Table 7.1:** Termwise number of coefficients  $p$  and effective degree of freedom  $edf$  for fitted test model

	$\hat{f}_0$	$\hat{f}_1$	$\hat{f}_2$	$\hat{f}_3$	$\hat{f}_4$	$\hat{f}_5$	$\hat{f}_6$	$\hat{f}_7$	$\hat{f}_8$	$\hat{f}_9$	$\hat{f}_{10}$
$p$	5	4	149	2874	39	2308	195	5811	18	2682	702
$edf$	3.97	1.00	93.38	2679.24	6.79	2214.06	119.86	3646.52	10.06	1413.80	26.22
	$\hat{f}_{11}$	$\hat{f}_{12}$	$\hat{f}_{13}$	$\hat{f}_{14}$	$\hat{f}_{15}$	$\hat{f}_{16}$	$\hat{f}_{17}$	$\hat{f}_{18}$	$\hat{f}_{19}$	$\hat{f}_{20}$	
$p$	9	9	196	1341	1341	2088	2088	351	351	9	
$edf$	6.18	2.88	73.65	498.77	513.09	655.69	822.12	215.46	191.36	5.92	

In fact, that  $\hat{\mathbf{y}}$  computation is time consuming also explains why pseudo QR reduction took longer time than what was predicted. The online updating algorithm for pseudo QR reduction (see Figure 4.6) not just computes matrix cross-product. A chunk of model matrix has to be formed first. Computation of  $\hat{\mathbf{y}}$  also proceeds by chunks, requiring chunkwise formation of the model matrix. The timing result suggests that this formation costs is somehow rather high. Subtracting 25 minutes from 1 hour 15 minutes gives 50 minutes, i.e., 3000 seconds. This is closer to my estimate of 2735 seconds for matrix cross-product.

From this test model, it is observed that formation of model matrix accounts for nearly 1/3 of the GAM fitting time. For increasingly large  $n$  this proportion is even higher. For example, if  $n$  is now 10 times as large and  $p$  unchanged, both matrix cross-product and model matrix formation will take 500 minutes, but REML estimation will still take 45 minutes, thus the proportion approximates 1/2. This is evidently where the GAM fitting method could / should be improved. I will come back to this issue in the next Chapter.

### 7.3 Checking the test model

The test model involves 22571 regression coefficients for 1342159 data, ending up with 13201 model degree of freedom. The estimated residual variance is 0.2797, and the variance of logBS data is 1.3945, so the fitted model has an adjusted R-squared of 79.92%. The estimated variance for station-specific random effect and week-specific random effect are respectively 0.0991 and 0.1069. Terwise number of coefficients and effective degree of freedom are summarized in Table 7.1.

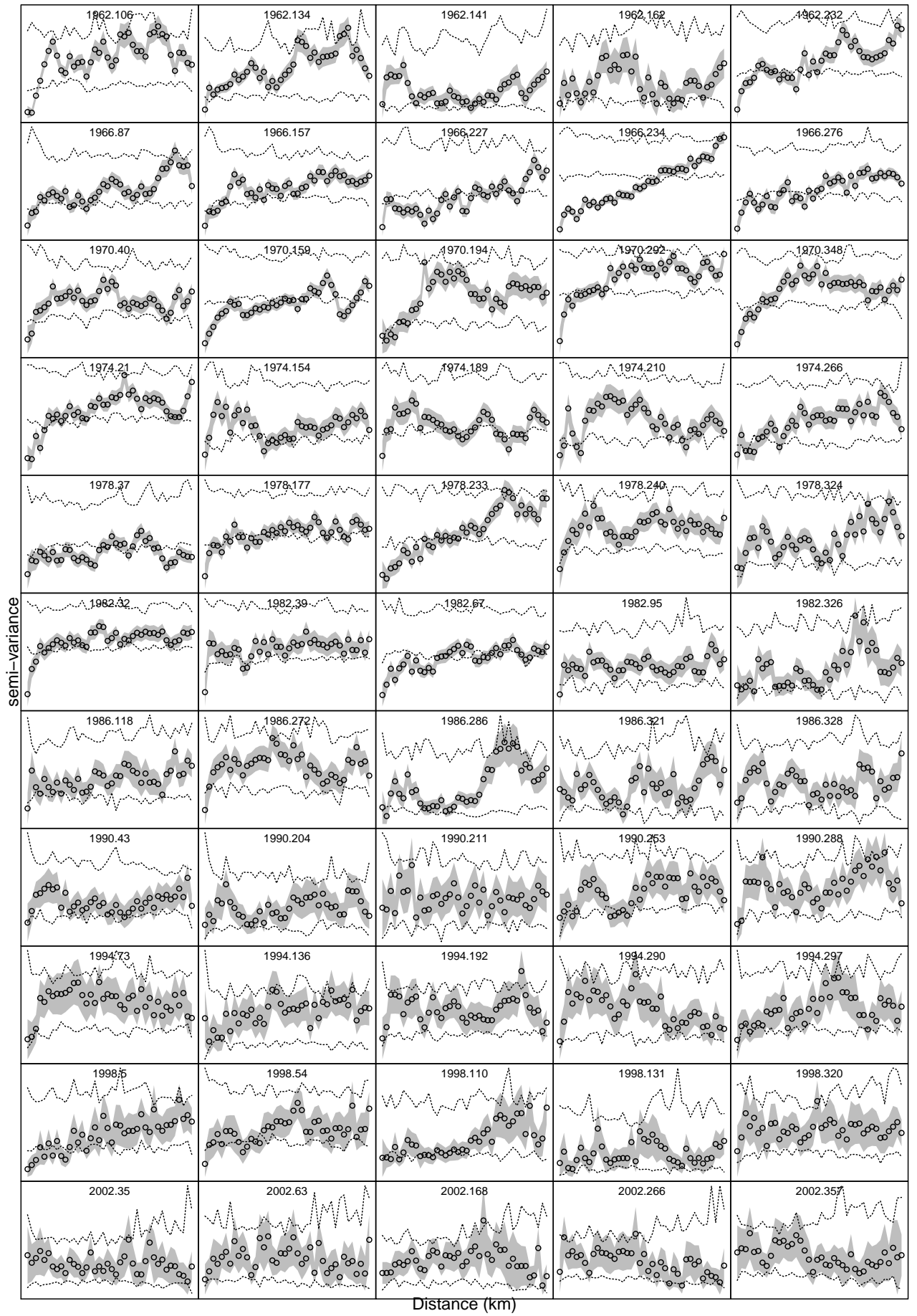
Before further checking everything else, I would first check whether there is spatial autocorrelation in residuals. Figure 7.1 presents a set of variograms of spatial residuals from 5 randomly chosen days from every 4 years. It is clear that the model is very inadequate in spatial modelling for the first two decades when the Network was relatively dense.

Figure 7.2 further checks random effects in the model. It turns out that the i.i.d. assumption for station-specific random effect is violated. This is a message that the model has difficulty in spatial modelling.

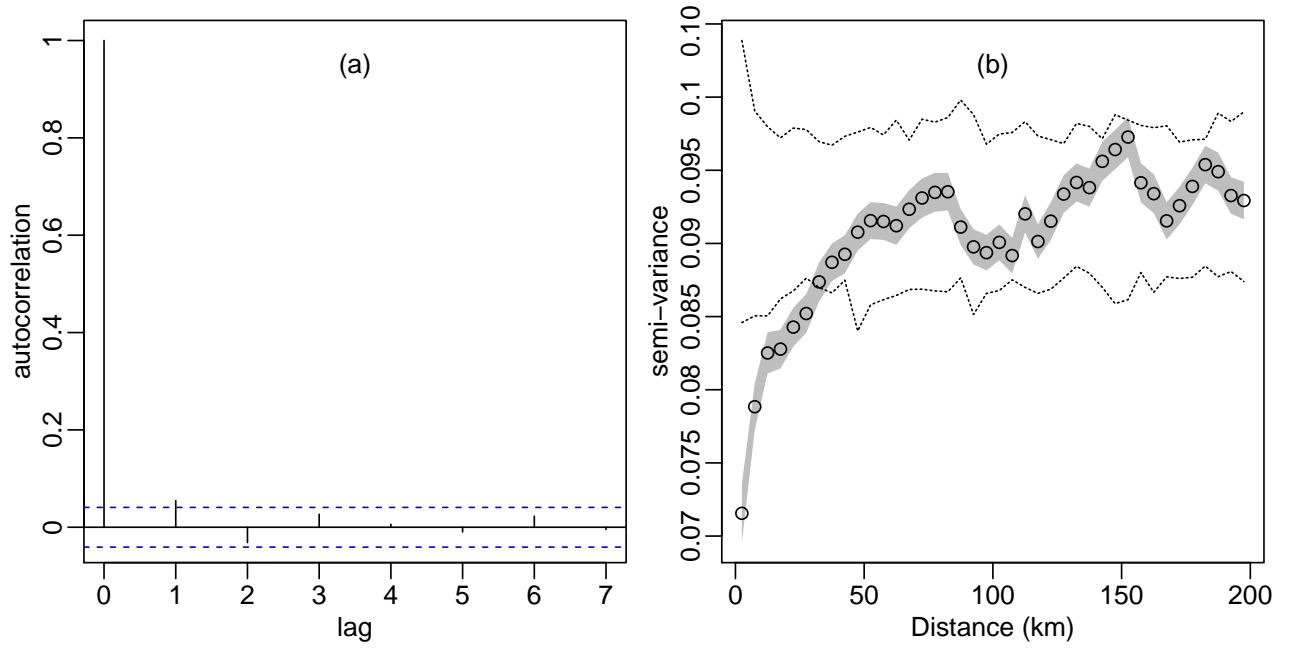
Figure 7.3 produces the transects plot for the fitted tensor product spline  $\hat{f}_7$  that models space-time interaction. The shape of those curve indicates of spatial-temporal extrapolation.

### 7.4 Summary

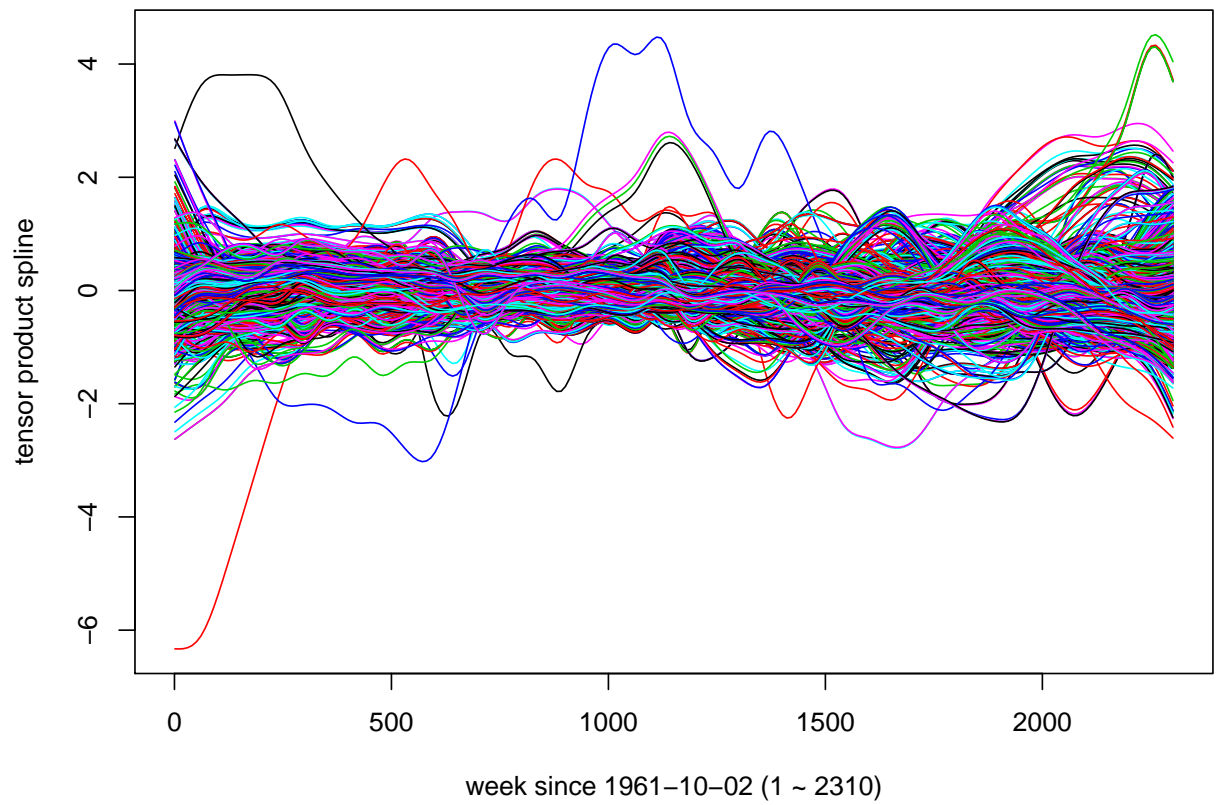
I stopped further checking the fitted model, because the message so far is readily clear: the model is very poor. I have to admit that I am unable to develop a reasonably good model. From now on, I will give up model development and focus on improving GAM fitting methods.



**Figure 7.1:** For the fitted model 7.1, a set of variograms of spatial residuals from 5 randomly chosen days from every 4 years.



**Figure 7.2:** Checking i.i.d. random effects in the the fitted model 7.1. Panel (a) is an ACF of week-specific random effects; no autocorrelation is seen. Panel (b) is a variogram of station-specific random effects; however, i.i.d. assumption is violated.



**Figure 7.3:** Transects plot for the fitted tensor product spline  $\hat{f}_7$ . These curves are added to the general long term trend for station-specific trend. The shape of those curves indicates there will be unreasonable spatial-temporal extrapolation.

## Chapter 8

# Discrete pseudo QR reduction

In §7.2 of the previous Chapter, I realized a performance bottleneck for GAM fitting with large datasets: chunkwise formation of  $\mathbf{X}$  is practically expensive compared with the computation of matrix cross-product  $\mathbf{X}'\mathbf{W}\mathbf{X}$  in “online” pseudo QR reduction; it is also what makes  $\hat{\mathbf{y}}$  computation very time consuming. Thus, faster construction of  $\mathbf{X}$  is key to faster pseudo QR reduction and  $\hat{\mathbf{y}}$  computation.

To address this issue, this Chapter will be organized in three stages as follows.

1. §8.1 explains how model matrix construction is conventionally conducted in `gam` function, how it is done in `bam` instead and how their difference eventually leads to the covariate discretization strategy for `bam` that is formally defined in §8.2;
2. §8.3, §8.4 and §8.5 present a set of algorithms for computations of  $\mathbf{X}'\mathbf{W}\mathbf{Y}$ ,  $\mathbf{X}'\mathbf{W}\mathbf{X}$  and  $\hat{\mathbf{y}}$  by exploiting packed storage for design matrices after covariate discretization. These algorithms to a large extent eliminate the need for constructing  $\mathbf{X}$ , thus is expected to be significantly faster. Additional attempts to enhance performance of these algorithms via block algorithms / data caching are covered in §8.6.
3. §8.7 experiments the new discrete computation method on the test logBS model 7.1.
4. §8.8 summarizes advantages and limitations of the new computation method, preluding the next Chapter, the most important Chapter of this thesis.

## 8.1 Faster model matrix formation for pseudo QR reduction

Suppose we have an additive model

$$\mathbf{y} = \mathbf{X}_0\boldsymbol{\beta}_0 + \sum_{i=1}^l \mathbf{f}_i + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{W}^{-1}\phi),$$

where in the linear predictor  $\boldsymbol{\eta} = \mathbf{X}_0\boldsymbol{\beta}_0 + \sum_{i=1}^l \mathbf{f}_i$ , there are

- $\mathbf{X}_0$ , a design matrix from some parametric terms;
- $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_l$ , smooth functions represented as splines.

Construction of the model matrix  $\mathbf{X}$  for this model is to form design matrix  $\mathbf{X}_i$  for all smooth functions.

Suppose we are to construct the spline basis for a univariate spline (like a cubic regression spline) given a vector of covariate values  $\mathbf{z} = (z_1, z_2, \dots, z_n)$ . An elementary fact is that if  $z_i = z_j$ , the  $i^{\text{th}}$  and  $j^{\text{th}}$  rows of the design matrix  $\mathbf{Z}$  are identical. In other words, construction of  $\mathbf{Z}$  only depends on the set of unique covariate values  $\bar{\mathbf{z}} = (\bar{z}_1, \bar{z}_2, \dots, \bar{z}_m)$ , where  $m$  is the number of unique values. Precisely, we can first construct a *packed design matrix* (with no redundant rows)  $\bar{\mathbf{Z}}$  from  $\bar{\mathbf{z}}$ , then *unpack* it. To this end, assume we have an index vector  $\mathbf{k} = (k_1, k_2, \dots, k_n)$  between  $\mathbf{z}$  and  $\bar{\mathbf{z}}$ , such that  $z_i = \bar{z}_{k_i}$ , then the  $i^{\text{th}}$  row of  $\mathbf{Z}$  is just the  $k_i^{\text{th}}$  row of  $\bar{\mathbf{Z}}$ . This construction idea works with other splines, too.

- For a multivariate spline but not a tensor product one, like a bivariate thin-plate regression spline, we have a covariate matrix instead of a covariate vector. The unique covariate values are then the unique rows of this matrix.
- For a tensor product spline, its marginal design matrices can be constructed first with the above method, then after unpacking we can form their row-wise Kronecker product.

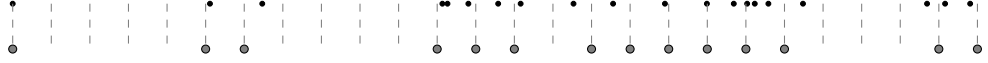
The model matrix construction procedure just described has always been implemented in `gam` function of `mgcv`, but not in `bam`. A fundamental assumption of `bam` is that  $n$  is very large, then  $m$  is also likely large, which is particularly true if some covariate is a continuous variable. Therefore, there may not even be enough RAM to store all packed design matrices. As a result, `bam` constructs  $\mathbf{X}$  with a different strategy. Let  $c$  be a row chunk size (the same one used in “online” model matrix reduction; see §4.2), it has two phases:

1. Prior to “online” model matrix reduction, randomly samples  $c$  covariate values from  $\mathbf{z}$ , then appending the minimum and maximum of  $\mathbf{z}$ , construct a design matrix using this subset (which does not necessarily contain all unique values of  $\mathbf{z}$ );
2. During the “online” reduction, construct  $\mathbf{X}$  chunk by chunk using the initial construction information (like knots locations and transformation matrices) from phase one.

While this keeps memory footprint under control, phase two treats every row of  $\mathbf{X}$  as if it is unique and computes it, ending up with a complexity proportional to  $n$  not  $m$ . Thus if  $m \ll n$ , it is very computationally inefficient.

Many covariates in the logBS dataset only take values from a discrete set. For example, there are only 2874 locations / stations; all time variables are discrete. Other covariates like elevation, temperature and rainfall, while conceptually being continuous variables, are up to an accuracy of 1 metre,  $0.1^\circ\text{C}$  and 0.1 millimetre, respectively. They are therefore also discrete with no more than a few hundreds to a few thousands of unique values. In the end, there is  $m \ll n$ . There will be sufficient RAM and it is more computationally efficient to first construct all packed design matrices, then unpack them during model matrix reduction. A quick experiment showed that for the test model 7.1, the formation time of  $\mathbf{X}$  was reduced from 25 minutes to just 3 minutes.

In general, many datasets with continuous covariates do not admit  $m \ll n$ . But is it possible to discretize these covariates, by approximating original covariate values with the nearest values on a regular grid, as if these covariates were measured up to certain accuracy? If this approximation is legitimate in general, the idea of packing and unpacking can be encapsulated into `bam` for efficient model matrix construction.



**Figure 8.1:** An illustration of covariate discretization by rounding. Black dots are original data. Gray dashed lines bin the original data onto a regular grid covering the range of the data. Gray dots are discretized data on the grid. In this example, original data  $z_i$ 's are 20 samples from Beta(2, 3) distribution, sorted in ascending order. The resolution of the grid is the 60% quantile of  $(z_{i+1} - z_i)$ .

## 8.2 Covariate discretization

In this section, I will formulate the idea of covariate discretization brought up at the end of the last section.

For a univariate continuous covariate  $\mathbf{z} = (z_1, z_2, \dots, z_n)$ , its discretization can be achieved by rounding. See Figure 8.1 for an illustration. In general, define that

- $m$  is the number of unique discretized covariate values, and  $\bar{\mathbf{z}} = (\bar{z}_1, \bar{z}_2, \dots, \bar{z}_m)$  is the set of unique values;
- $\mathbf{k} = (k_1, k_2, \dots, k_n)$  is discretization index matching  $\mathbf{z}$  and  $\bar{\mathbf{z}}$ , such that  $z_i \approx \bar{z}_{k_i}$ ;
- $\bar{\mathbf{Z}}$  and  $\mathbf{Z}$  are respectively the packed design matrix and full design matrix for a spline of  $\bar{\mathbf{z}}$  and  $\mathbf{z}$ .

For a multivariate continuous covariate, we have a covariate matrix. Discretization is first applied to each of its columns, then unique covariate values are determined by the unique rows of the discretized matrix. Finally, packed design matrix and discretization index can be similarly defined. Under discretization, the  $i^{\text{th}}$  row of  $\mathbf{Z}$  can be approximated by the  $k_i^{\text{th}}$  row of  $\bar{\mathbf{Z}}$ . Inserting some subscript, I will define the collection of  $(\bar{\mathbf{Z}}, \mathbf{k}_{\bar{\mathbf{Z}}}, m_{\bar{\mathbf{Z}}}, p_{\bar{\mathbf{Z}}})$  as the *packed storage* for  $\mathbf{Z}$ , where  $p_{\bar{\mathbf{Z}}}$  is the number of columns of  $\bar{\mathbf{Z}}$ .

Figure 8.2 is a toy example assessing the impact of covariate discretization on estimation. Estimation is still close to truth even if 1000 covariate values are discretized to just 25 unique values.

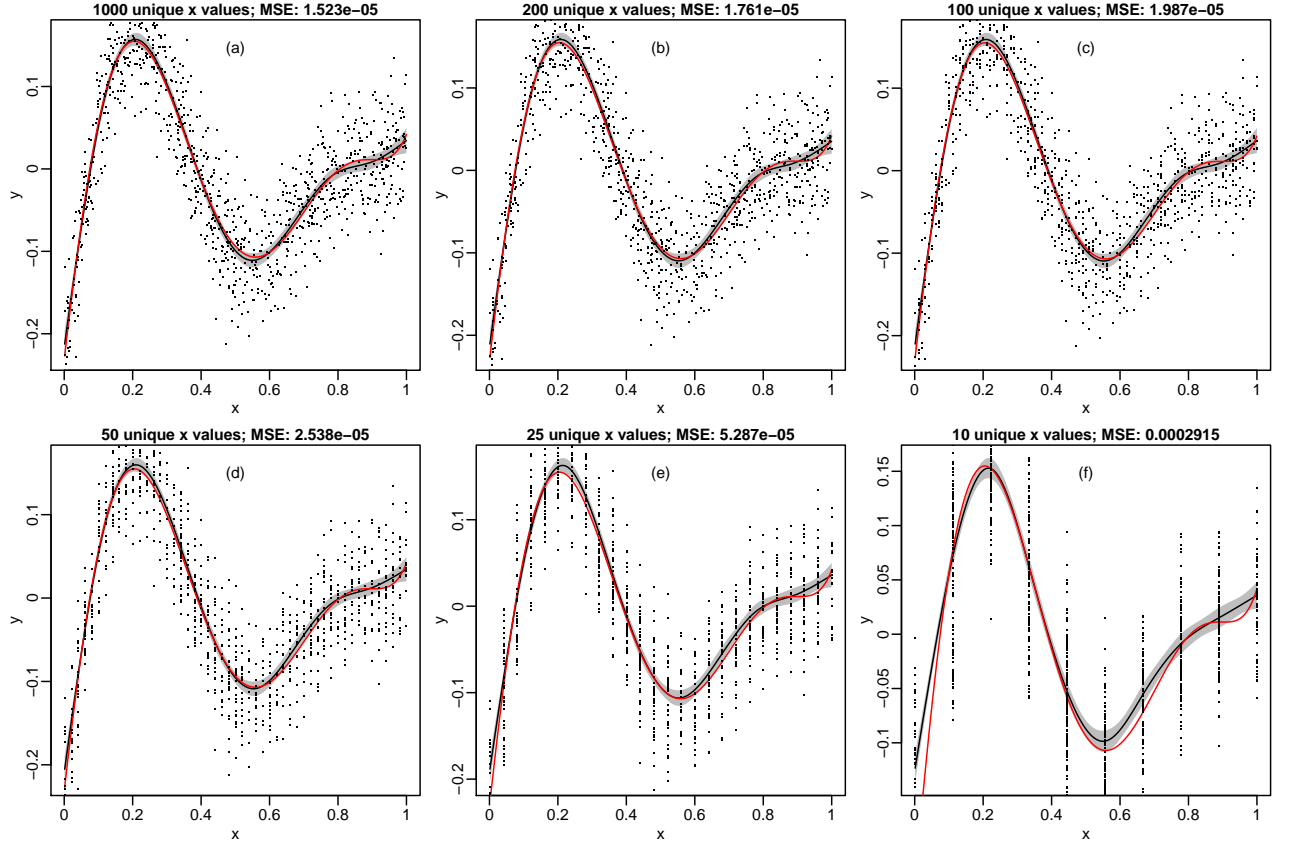
It should be emphasized that the “covariate discretization”, “packed design matrix”, “packed storage” and “unpacking” mentioned so far are only related to design matrix  $\mathbf{X}_i$  of a spline  $\mathbf{f}_i$ . The parametric design matrix  $\mathbf{X}_0$  is always explicitly constructed as an  $n \times p_0$  matrix.

## 8.3 Computation of $\mathbf{X}'\mathbf{W}\mathbf{y}$

Consider the linear predictor in §8.1, whose model matrix is  $\mathbf{X} = (\mathbf{X}_0 \mathbf{X}_1 \mathbf{X}_2 \dots \mathbf{X}_l)$ . Computation of the matrix-vector cross-product  $\mathbf{X}'\mathbf{W}\mathbf{y}$ , or  $\mathbf{X}'\tilde{\mathbf{y}}$  where  $\tilde{\mathbf{y}} = \mathbf{W}\mathbf{y}$  is the weighted response vector, may proceed block by block

$$\mathbf{X}'\tilde{\mathbf{y}} = \begin{pmatrix} \mathbf{X}_0'\tilde{\mathbf{y}} \\ \mathbf{X}_1'\tilde{\mathbf{y}} \\ \vdots \\ \mathbf{X}_l'\tilde{\mathbf{y}} \end{pmatrix}.$$

Given packed storage for  $\mathbf{X}_i$ 's, conventional computation first requires unpacking all  $\mathbf{X}_i$ 's (if  $\mathbf{X}_i$  is a tensor product design matrix, first unpack all its marginal design matrices, then form their row-wise Kronecker product). In this section, I will describe new computation method for  $\mathbf{X}_i'\tilde{\mathbf{y}}$ , where the initial matrix unpacking can be avoided.



**Figure 8.2:** A toy example assessing the impact of covariate discretization on estimation. In panel (a), for 1000 uniformly distributed  $x$  values on  $[0, 1]$ , a random polynomial of degree 6 is generated as the true function  $y(x)$  (red curve). Gaussian white noise is added with a noise-to-signal ratio of 0.5 (black dots are noise observations). A natural cubic regression spline of 10 knots is fitted (black curve), and the gray ribbon gives 95% pointwise confidence interval of the estimation. In panel (b) to (e),  $x$  are discretized to 200, 100, 50, 25 and 10 unique values. The effect of discretization is most obvious from the scatter plot, as they are gradually shrunk toward vertical bars. However, apart from the last panel where 10 unique values are too few, the estimated spline function (black curve) still well approximates the true function (red curve).

For convenience, I will call  $\mathbf{X}_i$  a *singleton*, if it is a single design matrix, and a *tensor*, if it is a tensor product design matrix with at least two margins. Computation of  $\mathbf{X}_0' \tilde{\mathbf{y}}$  needs no special treatment, as parametric design matrix  $\mathbf{X}_0$  is not stored in packed format. The focus is to demonstrate algorithms for computing cross-product between a singleton  $\mathbf{X}_i$  and  $\tilde{\mathbf{y}}$ , and that between a tensor  $\mathbf{X}_i$  and  $\tilde{\mathbf{y}}$ .

Assume that a singleton  $\mathbf{X}_i$  has packed storage  $(\bar{\mathbf{A}}, \mathbf{k}_{\bar{\mathbf{A}}}, m_{\bar{\mathbf{A}}}, p_{\bar{\mathbf{A}}})$ . An algorithm for  $\mathbf{X}_i' \tilde{\mathbf{y}}$  is as straightforward as Figure 8.3a. For a tensor  $\mathbf{X}_i$  with  $s$  packed margins from  $(\bar{\mathbf{A}}_1, \mathbf{k}_{\bar{\mathbf{A}}_1}, m_{\bar{\mathbf{A}}_1}, p_{\bar{\mathbf{A}}_1})$  to  $(\bar{\mathbf{A}}_s, \mathbf{k}_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{A}}_s}, p_{\bar{\mathbf{A}}_s})$ , combine the first  $(s-1)$  margins into a single matrix  $\dot{\mathbf{A}} = \mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2 \tilde{\otimes} \cdots \tilde{\otimes} \mathbf{A}_{s-1}$ , where  $\tilde{\otimes}$  is the row-wise Kronecker product between two matrices. Then  $\mathbf{X}_i = \dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s$  can be treated as a tensor with only two margins. Applying (1.1) to  $\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s$  would give

$$(\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \tilde{\mathbf{y}} = \begin{bmatrix} \mathbf{A}_s' \text{diag}(\dot{\mathbf{A}}(, 1)) \\ \mathbf{A}_s' \text{diag}(\dot{\mathbf{A}}(, 2)) \\ \vdots \\ \mathbf{A}_s' \text{diag}(\dot{\mathbf{A}}(, p_{\dot{\mathbf{A}}})) \end{bmatrix} \tilde{\mathbf{y}} = \begin{bmatrix} \mathbf{A}_s' (\text{diag}(\dot{\mathbf{A}}(, 1)) \tilde{\mathbf{y}}) \\ \mathbf{A}_s' (\text{diag}(\dot{\mathbf{A}}(, 2)) \tilde{\mathbf{y}}) \\ \vdots \\ \mathbf{A}_s' (\text{diag}(\dot{\mathbf{A}}(, p_{\dot{\mathbf{A}}})) \tilde{\mathbf{y}}) \end{bmatrix},$$

where  $p_{\dot{\mathbf{A}}} = \prod_{j=1}^{s-1} p_{\bar{\mathbf{A}}_j}$ . This implies that  $(\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \tilde{\mathbf{y}}$  can be obtained by computing  $\mathbf{A}_s' \tilde{\mathbf{y}}_k$  block by block, where  $\tilde{\mathbf{y}}_k = \text{diag}(\dot{\mathbf{A}}(, k)) \tilde{\mathbf{y}}$ . Figure 8.3b demonstrates the algorithm.

Implementation of algorithm 8.3b requires extracting the  $k^{\text{th}}$  column from the tensor product design matrix  $\dot{\mathbf{A}}$ . It can be shown that this column is an Hadamard product<sup>1</sup>

$$\dot{\mathbf{A}}(, k) = \mathbf{A}_1(, k_1) \circ \mathbf{A}_2(, k_2) \circ \cdots \circ \mathbf{A}_{s-1}(, k_{s-1}),$$

<sup>1</sup>For two vectors  $\mathbf{x} = (x_1 \ x_2 \ \cdots \ x_n)$  and  $\mathbf{y} = (y_1 \ y_2 \ \cdots \ y_n)$ , their Hadamard product is the element-wise



for some  $k_1, k_2, \dots, k_{s-1}$ . Here, using zero-based index (that is, if a matrix has  $p$  columns, they are columns 0, 1,  $\dots$ ,  $(p-1)$ ) turns out convenient. The  $k_j$  index can be determined most easily in a two-margin scenario  $\tilde{\mathbf{A}}(, k) = \mathbf{A}_1(, k_1) \circ \mathbf{A}_2(, k_2)$ . By applying (1.1) to  $\mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2$ , it can be derived that  $k_2 = \lfloor k/p_{\mathbf{A}_2} \rfloor$ ,  $k_1 = \lfloor k/p_{\mathbf{A}_1} \rfloor$ , where  $\lfloor a/b \rfloor$  is the *modulo* function and  $\lfloor a/b \rfloor$  is the *floor* function. For more than two margins, say three margins, a recursion applies. Collect the first two margins into a single matrix. Then if  $\tilde{\mathbf{A}}(, k) = (\mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2)(, \tilde{k}) \circ \mathbf{A}_3(, k_3)$ , the previous two-margin result gives  $k_3 = \lfloor k/p_{\mathbf{A}_3} \rfloor$  and  $\tilde{k} = \lfloor k/p_{\tilde{\mathbf{A}}} \rfloor$ . Now applying the same to  $(\mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2)(, \tilde{k}) = \mathbf{A}_1(, k_1) \circ \mathbf{A}_2(, k_2)$  gives  $k_2 = \lfloor \tilde{k}/p_{\mathbf{A}_2} \rfloor$  and  $k_1 = \lfloor \tilde{k}/p_{\mathbf{A}_1} \rfloor$ . In general, the  $k^{\text{th}}$  column can be computed using algorithm in Figure 8.3c.

Algorithms in Figures 8.3a and 8.3b respectively involve  $O(n) + O(m_{\tilde{\mathbf{A}}} p_{\tilde{\mathbf{A}}})$  and  $O(s n p_{\tilde{\mathbf{A}}}) + O(m_{\tilde{\mathbf{A}}} p_{\tilde{\mathbf{A}}} p_{\tilde{\mathbf{A}}_s})$  FLOP, smaller than the  $O(n p_{\tilde{\mathbf{A}}})$  and  $O(n p_{\tilde{\mathbf{A}}} p_{\tilde{\mathbf{A}}_s})$  FLOP via conventional computation method with initial matrix unpacking. These algorithms are expected to be  $O(\min\{p_{\tilde{\mathbf{A}}}, n/m_{\tilde{\mathbf{A}}}\})$  times faster for a singleton  $\mathbf{X}_i$ , and  $O(\min\{p_{\tilde{\mathbf{A}}_s}, n/m_{\tilde{\mathbf{A}}_s}\})$  times faster for a tensor  $\mathbf{X}_i$ . The latter result implies that computations of  $(\mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2)' \tilde{\mathbf{y}}$  and  $(\mathbf{A}_2 \tilde{\otimes} \mathbf{A}_1)' \tilde{\mathbf{y}}$  are not equally fast if  $p_{\tilde{\mathbf{A}}_1} \neq p_{\tilde{\mathbf{A}}_2}$ . Generally for a tensor  $\mathbf{X}_i$ , if its largest margin (the one with the maximal  $p_{\tilde{\mathbf{A}}_j}$ ) is the last margin, the maximal speedup can be attained. Therefore in `mgcv` implementation, it is worth exchanging the largest margin to the last margin during construction of a tensor product spline.

The centring constraint described in §4.1.6 implies that in practice, `mgcv` would estimate  $\tilde{\beta}_i$  not  $\beta_i$ , so the cross-product that is really needed is  $(\mathbf{X}_i \mathbf{Q}_i^\perp)' \tilde{\mathbf{y}}$ . It turns out that we can first work out  $\mathbf{z} = \mathbf{X}_i' \tilde{\mathbf{y}}$  using previously described algorithms, then compute  $(\mathbf{Q}_i^\perp)' \mathbf{z}$  as a post-processing.

## 8.4 Computation of $\mathbf{X}' \mathbf{W} \mathbf{X}$

Consider the linear predictor in §8.1, whose model matrix is  $\mathbf{X} = (\mathbf{X}_0 \mathbf{X}_1 \mathbf{X}_2 \dots \mathbf{X}_l)$ . Computation of the matrix cross-product  $\mathbf{X}' \mathbf{W} \mathbf{X}$  may proceed block by block

$$\mathbf{X}' \mathbf{W} \mathbf{X} = \begin{pmatrix} \mathbf{X}_0' \mathbf{W} \mathbf{X}_0 & \mathbf{X}_0' \mathbf{W} \mathbf{X}_1 & \dots & \mathbf{X}_0' \mathbf{W} \mathbf{X}_l \\ \mathbf{X}_1' \mathbf{W} \mathbf{X}_0 & \mathbf{X}_1' \mathbf{W} \mathbf{X}_1 & \dots & \mathbf{X}_1' \mathbf{W} \mathbf{X}_l \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}_l' \mathbf{W} \mathbf{X}_0 & \mathbf{X}_l' \mathbf{W} \mathbf{X}_1 & \dots & \mathbf{X}_l' \mathbf{W} \mathbf{X}_l \end{pmatrix},$$

where only the lower triangular part needs be computed due to symmetry. Given packed storage for  $\mathbf{X}_i$ 's, conventional computation requires unpacking all  $\mathbf{X}_i$ 's (if  $\mathbf{X}_i$  is a tensor product design matrix, first unpack all its marginal design matrices, then form their row-wise Kronecker product) before computing  $\mathbf{X}_i' \mathbf{W} \mathbf{X}_j$ . However, using algorithms in §8.3, computation of  $\mathbf{X}_i' \mathbf{W} \mathbf{X}_j$  can be faster. Consider only unpacking  $\mathbf{X}_j$  into a full design matrix  $\mathbf{Y}$ , then  $\mathbf{X}_i' \mathbf{W} \mathbf{X}_j = \mathbf{X}_i' \mathbf{W} \mathbf{Y}$  can be computed using algorithms in Figures 8.4a and 8.4b, which are straightforward generalization of algorithms in Figures 8.3a and 8.3b. Compared with conventional computation where both  $\mathbf{X}_i$  and  $\mathbf{X}_j$  are unpacked, these algorithms do not require unpacking  $\mathbf{X}_i$ , and are expected to be  $O(\min\{p_{\tilde{\mathbf{A}}}, n/m_{\tilde{\mathbf{A}}}\})$  times faster for a singleton  $\mathbf{X}_i$ , and  $O(\min\{p_{\tilde{\mathbf{A}}_s}, n/m_{\tilde{\mathbf{A}}_s}\})$  times faster for a tensor  $\mathbf{X}_i$ .

In general, fully unpacking  $\mathbf{X}_j$  would require too much RAM if  $n$  is large. A workaround is to only unpack a column chunk at one time, as implemented by algorithms in Figure 8.5a and 8.5b. With a chunk size  $b$ , the memory footprint is bounded by  $O(nb)$ . Note that for a tensor  $\mathbf{X}_i$ , the value of  $b$  would affect the complexity of the algorithm<sup>2</sup>. As a result, the algorithm is  $O(\min\{p_{\tilde{\mathbf{A}}_s}/(1 + s/b), n/m_{\tilde{\mathbf{A}}_s}\})$  times faster than conventional computation. To maximize this speedup factor,  $s/b$  should be as small as possible. A reasonable choice is  $b = 20s$  so that  $s/b = 0.05$ .

product  $\mathbf{x} \circ \mathbf{y} = (x_1 y_1 \quad x_2 y_2 \quad \dots \quad x_n y_n)$ . Hadamard product is commutative and associative, and it is straightforward to extend to more than two vectors.

<sup>2</sup>Because after chunking, every column of  $\tilde{\mathbf{A}}$  is respectively extracted  $(p_{\tilde{\mathbf{B}}}/b)$  and  $(p_{\tilde{\mathbf{B}}}/b)$  times rather than only once for a singleton  $\mathbf{X}_j$  and a tensor  $\mathbf{X}_j$ .

---

```

initialize a length- $m_{\bar{A}}$  vector of zeros
 $\tilde{\mathbf{y}} = \text{zeros}(m_{\bar{A}})$ 
vector aggregation,  $O(n)$ 
for  $i = 1 : n$ 
     $\tilde{\mathbf{y}}(\mathbf{k}_{\bar{A}}(i)) += \tilde{\mathbf{y}}(i)$ 
matrix-vector cross-product,  $O(m_{\bar{A}}p_{\bar{A}})$ 
 $\mathbf{A}'\tilde{\mathbf{y}} = \bar{\mathbf{A}}'\tilde{\mathbf{y}}$ 

```

---

(a)  $\mathbf{X}'_i\tilde{\mathbf{y}} = \mathbf{A}'\tilde{\mathbf{y}}$ , where  $\mathbf{X}_i$  is a singleton with packed storage  $(\bar{\mathbf{A}}, \mathbf{k}_{\bar{A}}, m_{\bar{A}}, p_{\bar{A}})$ . FLOP count for each computation step is given in gray text, and the algorithm involves  $O(n) + O(m_{\bar{A}}p_{\bar{A}})$  FLOP.

---

```

 $p_{\bar{A}} = \prod_{j=1}^{s-1} p_{\bar{A}_j}$ 
for  $k = 1 : p_{\bar{A}}$ 
    compute  $\tilde{\mathbf{y}}_k$ ,  $O(sn)$ 
    compute  $\tilde{\mathbf{y}}_k = \text{diag}(\dot{\mathbf{A}}(:, k))\tilde{\mathbf{y}}$  using algorithm in Figure 8.3c
    compute elements  $(kp_{\bar{A}_s} - p_{\bar{A}_s} + 1)$  to  $kp_{\bar{A}_s}$  of  $\mathbf{X}'_i\tilde{\mathbf{y}}_k$ ,  $O(n) + O(m_{\bar{A}_s}p_{\bar{A}_s})$ 
    compute  $\mathbf{z} = \mathbf{A}'_s\tilde{\mathbf{y}}_k$  using algorithm in Figure 8.3a
     $(\mathbf{X}'_i\tilde{\mathbf{y}})((kp_{\bar{A}_s} - p_{\bar{A}_s} + 1) : kp_{\bar{A}_s}) = \mathbf{z}$ 

```

---

(b)  $\mathbf{X}'_i\tilde{\mathbf{y}} = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)'\tilde{\mathbf{y}}$ , where  $\mathbf{X}_i$  is a tensor with  $s$  packed margins  $(\bar{\mathbf{A}}_j, \mathbf{k}_{\bar{\mathbf{A}}_j}, m_{\bar{\mathbf{A}}_j}, p_{\bar{\mathbf{A}}_j})$  and  $\dot{\mathbf{A}} = \mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2 \tilde{\otimes} \cdots \tilde{\otimes} \mathbf{A}_{s-1}$ . FLOP count for each computation step is given in gray text, and the algorithm involves  $O(snp_{\bar{A}}) + O(m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}}p_{\bar{\mathbf{A}}_s})$  FLOP.

---

```

initialize a length- $n$  vector of ones,  $O(n)$ 
 $\mathbf{a} = \text{ones}(n)$ 
initialize  $\tilde{k}$ 
 $\tilde{k} = k$ 
for  $j = (s - 1) : 1$ 
    compute  $k_j$ 
     $k_j = \lfloor \tilde{k} / p_{\bar{\mathbf{A}}_j} \rfloor$ 
    update  $\tilde{k}$ 
     $\tilde{k} = \lfloor \tilde{k} / p_{\bar{\mathbf{A}}_j} \rfloor$ 
    compute Hadamard product,  $O(n)$ 
     $\mathbf{a} = \mathbf{a} \circ \bar{\mathbf{A}}_j(\mathbf{k}_{\bar{\mathbf{A}}_j}, k_j)$ 
store result
 $\dot{\mathbf{A}}(:, k) = \mathbf{a}$ 

```

---

(c) Compute the  $k^{\text{th}}$  column of  $\dot{\mathbf{A}} = \mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2 \tilde{\otimes} \cdots \tilde{\otimes} \mathbf{A}_{s-1}$ . Here, zero-based index is used for convenience, that is, if a matrix has  $p$  columns, they are columns 0, 1,  $\dots$ ,  $(p - 1)$ . When  $\bar{\mathbf{A}}_j$  is stored in packed format  $(\bar{\mathbf{A}}_j, \mathbf{k}_{\bar{\mathbf{A}}_j}, m_{\bar{\mathbf{A}}_j}, p_{\bar{\mathbf{A}}_j})$ , the Hadamard product can be computed as  $\mathbf{a} = \mathbf{a} \circ \bar{\mathbf{A}}_j(\mathbf{k}_{\bar{\mathbf{A}}_j}, k_j)$ . Note that this algorithm essentially computes  $\dot{\mathbf{A}}(:, k) \circ \mathbf{a}$ , or equivalently  $\text{diag}(\dot{\mathbf{A}}(:, k))\mathbf{a}$ , but it returns  $\dot{\mathbf{A}}(:, k)$  if  $\mathbf{a}$  is initially a vector of ones. The algorithm involves  $O(sn)$  FLOP.

**Figure 8.3:** Computation of  $\mathbf{X}'_i\mathbf{W}\mathbf{y} = \mathbf{X}'_i\tilde{\mathbf{y}}$ , where  $\mathbf{X}_i$  is a singleton or tensor, and  $\tilde{\mathbf{y}} = \mathbf{W}\mathbf{y}$  is the weighted response vector.

---

initialize a  $m_{\bar{A}} \times p_Y$  matrix of zeros

$\tilde{Y} = \text{zeros}(m_{\bar{A}}, p_Y)$

matrix row-aggregation,  $O(np_Y)$

for  $i = 1 : n$

$\tilde{Y}(\mathbf{k}_{\bar{A}}(i), :) += \tilde{Y}(i, :)$

matrix cross-product,  $O(m_{\bar{A}}p_{\bar{A}}p_Y)$

$\mathbf{A}'\tilde{Y} = \bar{\mathbf{A}}'\tilde{Y}$

---

(a)  $\mathbf{X}'_i\tilde{Y} = \mathbf{A}'\tilde{Y}$ , where  $\mathbf{X}_i$  is a singleton with packed storage  $(\bar{\mathbf{A}}, \mathbf{k}_{\bar{A}}, m_{\bar{A}}, p_{\bar{A}})$ . FLOP count for each computation step is given in gray text, and the algorithm involves  $O(np_Y) + O(m_{\bar{A}}p_{\bar{A}}p_Y)$  FLOP.

---

$p_{\dot{\mathbf{A}}} = \prod_{j=1}^{s-1} p_{\bar{A}_j}$

for  $k = 1 : p_{\dot{\mathbf{A}}}$

form  $k^{\text{th}}$  column of  $\dot{\mathbf{A}} = \mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2 \tilde{\otimes} \cdots \tilde{\otimes} \mathbf{A}_{s-1}$ ,  $O(sn)$

compute  $\mathbf{a} = \dot{\mathbf{A}}(:, k)$  using algorithm in Figure 8.3c

matrix row-scaling,  $O(np_Y)$

$\tilde{Y}_k = \text{diag}(\mathbf{a})\tilde{Y}$

compute rows  $(kp_{\bar{A}_s} - p_{\bar{A}_s} + 1)$  to  $kp_{\bar{A}_s}$  of  $\mathbf{X}'_i\tilde{Y}$ ,  $O(np_Y) + O(m_{\bar{A}_s}p_{\bar{A}_s}p_Y)$

compute  $\mathbf{Z} = \mathbf{A}'_s\tilde{Y}_k$  using algorithm in Figure 8.4a

$(\mathbf{X}'_i\tilde{Y})((kp_{\bar{A}_s} - p_{\bar{A}_s} + 1) : kp_{\bar{A}_s}, :) = \mathbf{Z}$

---

(b)  $\mathbf{X}'_i\tilde{Y} = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)'\tilde{Y}$ , where  $\mathbf{X}_i$  is a tensor with  $s$  packed margins from  $(\bar{\mathbf{A}}_1, \mathbf{k}_{\bar{\mathbf{A}}_1}, m_{\bar{\mathbf{A}}_1}, p_{\bar{\mathbf{A}}_1})$  to  $(\bar{\mathbf{A}}_s, \mathbf{k}_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{A}}_s}, p_{\bar{\mathbf{A}}_s})$ . FLOP count for each computation step is given in gray text, and the algorithm involves  $O(sn p_{\dot{\mathbf{A}}}) + O(np_{\dot{\mathbf{A}}}p_Y) + O(m_{\bar{\mathbf{A}}_s}p_{\dot{\mathbf{A}}}p_Y)$  FLOP.

---

**Figure 8.4:** Computation of  $\mathbf{X}'_i\mathbf{W}\mathbf{Y} = \mathbf{X}'_i\tilde{Y}$ , where  $\mathbf{X}_i$  is a singleton or tensor and  $\tilde{Y} = \mathbf{W}\mathbf{Y}$  has  $p_Y$  columns.

Note that computations of  $\mathbf{X}'_i\mathbf{W}\mathbf{X}_j$  and  $\mathbf{X}'_j\mathbf{W}\mathbf{X}_i$  are not equally fast using algorithms above. For example, assuming that  $\mathbf{X}_i$  and  $\mathbf{X}_j$  are both tensors, the former yields  $O(p_{\bar{\mathbf{A}}_s})$  speedup, while the latter yields  $O(p_{\bar{\mathbf{B}}_s})$  speedup. This implies that if  $p_{\bar{\mathbf{A}}_s} < p_{\bar{\mathbf{B}}_s}$ , it is more efficient to compute  $\mathbf{X}'_j\mathbf{W}\mathbf{X}_i$ , then transpose it to get  $\mathbf{X}'_i\mathbf{W}\mathbf{X}_j$ . In `mgcv`, a run-time decision can be made on this.

With centring constraint, the cross-product needed is  $(\mathbf{X}_i\mathbf{Q}_i^\perp)'\mathbf{W}\mathbf{X}_j\mathbf{Q}_j^\perp$ . We can first get  $\mathbf{Z} = \mathbf{X}'_i\mathbf{W}\mathbf{X}_j$  using previously described algorithms, then compute  $(\mathbf{Q}_i^\perp)'\mathbf{Z}\mathbf{Q}_j^\perp$  as a post-processing.

## 8.5 Computation of $\hat{\mathbf{y}}$

Now let us consider the computation of  $\hat{\mathbf{y}}$ . For the linear prediction in §8.1, there is  $\hat{\mathbf{y}} = \sum_{i=0}^l \hat{\mathbf{y}}_i = \sum_{i=0}^l \mathbf{X}_i\hat{\boldsymbol{\beta}}_i$ . Conventional computation would first unpack all  $\mathbf{X}_i$ 's from their packed storage then do matrix-vector multiplication. But a more efficient method is possible without matrix unpacking. To simplify notations in the following, let  $\mathbf{b}_i = \hat{\boldsymbol{\beta}}_i$ .

If  $\mathbf{X}_i$  is a singleton, then  $\hat{\mathbf{y}}_i$  can be straightforwardly obtained with the algorithm in Figure 8.6a which is  $O(\min\{p_{\bar{\mathbf{A}}}, n/m_{\bar{\mathbf{A}}}\})$  times faster than conventional computation. If  $\mathbf{X}_i$  is a tensor with  $s$  margins from  $(\bar{\mathbf{A}}_1, \mathbf{k}_{\bar{\mathbf{A}}_1}, m_{\bar{\mathbf{A}}_1}, p_{\bar{\mathbf{A}}_1})$  to  $(\bar{\mathbf{A}}_s, \mathbf{k}_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{A}}_s}, p_{\bar{\mathbf{A}}_s})$ , there is

$$\hat{\mathbf{y}}_i = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)\mathbf{b}_i = \left( \cdots \quad \text{diag}(\dot{\mathbf{A}}(:, k))\mathbf{A}_s \quad \cdots \right) \begin{pmatrix} \vdots \\ \mathbf{b}_i((kp_{\bar{\mathbf{A}}_s} - p_{\bar{\mathbf{A}}_s} + 1) : kp_{\bar{\mathbf{A}}_s}) \\ \vdots \end{pmatrix}$$

$$= \sum_{k=1}^{p_{\dot{\mathbf{A}}}} \text{diag}(\dot{\mathbf{A}}(:, k))\mathbf{A}_s\mathbf{b}_i((kp_{\bar{\mathbf{A}}_s} - p_{\bar{\mathbf{A}}_s} + 1) : kp_{\bar{\mathbf{A}}_s})$$

---

```

 $k = 1$ 
while ( $k \leq p_B$ )
  determine column chunk size
   $\tilde{b} = \min\{b, p_B - k + 1\}$ 
  unpack columns  $k$  to  $(k + \tilde{b} - 1)$  of  $B$ ,  $O(n\tilde{b})$ 
  for  $\tilde{k} = 1 : \tilde{b}$ 
     $Y(:, \tilde{k}) = \bar{B}(\mathbf{k}_{\bar{B}}, k + \tilde{k} - 1)$ 
  matrix weighting ( $Y$  can be overwritten),  $O(n\tilde{b})$ 
   $Y = WY$ 
  compute columns  $k$  to  $(k + \tilde{b} - 1)$  of  $X_i' W X_j$ 
   $O(n\tilde{b}) + O(m_{\bar{A}} p_{\bar{A}} \tilde{b})$  or  $O(snp_A) + O(np_A \tilde{b}) + O(m_{\bar{A}_s} p_{\bar{A}} p_{\bar{A}_s} \tilde{b})$ 
  compute  $Z = X_i' Y$  using algorithm in Figure 8.4a or 8.4b
   $(X_i' W X_j)(, k : (k + \tilde{b} - 1)) = Z$ 
  proceed to column chunk
   $k = k + \tilde{b}$ 

```

---

(a)  $X_i' W X_j = X_i' W B$ , where  $X_j$  is a singleton with packed storage  $(\bar{B}, \mathbf{k}_{\bar{B}}, m_{\bar{B}}, p_{\bar{B}})$ . FLOP count for each computation step is given in gray text, and the algorithm involves  $O(np_B) + O(m_{\bar{A}} p_{\bar{A}} p_B)$  FLOP for a singleton  $X_i$  and  $O(np_B) + O(snp_A p_B/b) + O(np_A p_B) + O(m_{\bar{A}_s} p_{\bar{A}} p_{\bar{A}_s} p_B)$  FLOP for a tensor  $X_i$ .

---

```

 $p_B = \prod_{j=1}^t p_{\bar{B}_j}$ 
 $k = 1$ 
while ( $k \leq p_B$ )
  determine column chunk size
   $\tilde{b} = \min\{b, p_B - k + 1\}$ 
  unpack columns  $k$  to  $(k + \tilde{b} - 1)$  of  $B_1 \tilde{\otimes} B_2 \tilde{\otimes} \dots \tilde{\otimes} B_t$ ,  $O(n\tilde{b})$ 
  for  $\tilde{k} = 1 : \tilde{b}$ 
     $Y(:, \tilde{k}) = (B_1 \tilde{\otimes} B_2 \tilde{\otimes} \dots \tilde{\otimes} B_t)(, k + \tilde{k} - 1)$ 
  matrix weighting ( $Y$  can be overwritten),  $O(n\tilde{b})$ 
   $Y = WY$ 
  compute columns  $k$  to  $(k + \tilde{b} - 1)$  of  $X_i' W X_j$ 
   $O(n\tilde{b}) + O(m_{\bar{A}} p_{\bar{A}} \tilde{b})$  or  $O(snp_A) + O(np_A \tilde{b}) + O(m_{\bar{A}_s} p_{\bar{A}} p_{\bar{A}_s} \tilde{b})$ 
  compute  $Z = X_i' Y$  using algorithm in Figure 8.4a or 8.4b
   $(X_i' W X_j)(, k : (k + \tilde{b} - 1)) = Z$ 
  proceed to column chunk
   $k = k + \tilde{b}$ 

```

---

(b)  $X_i' W X_j = X_i' W (B_1 \tilde{\otimes} B_2 \tilde{\otimes} \dots \tilde{\otimes} B_t)$ , where  $X_j$  is a tensor with  $t$  packed margins from  $(\bar{B}_1, \mathbf{k}_{\bar{B}_1}, m_{\bar{B}_1}, p_{\bar{B}_1})$  to  $(\bar{B}_t, \mathbf{k}_{\bar{B}_t}, m_{\bar{B}_t}, p_{\bar{B}_t})$ . FLOP count for each computation step is given in gray text, and the algorithm involves  $O(n p_B) + O(m_{\bar{A}} p_{\bar{A}} p_B)$  FLOP for a singleton  $X_i$  and  $O(n p_B) + O(snp_A p_B/b) + O(np_A p_B) + O(m_{\bar{A}_s} p_{\bar{A}} p_{\bar{A}_s} p_B)$  FLOP for a tensor  $X_i$ .

---

**Figure 8.5:** Computation of  $X_i' W X_j$ , where  $X_i$  and  $X_j$  can be singletons or tensors. (If  $X_j$  is the parametric design matrix  $X_0$ , computation of  $X_i' W X_0$  is straightforward using algorithms in Figure 8.4a or 8.4b by setting  $Y = X_0$ .)

---

initialize a length- $n$  vector of zeros,  $O(n)$

$\hat{\mathbf{y}}_i = \text{zeros}(n)$

matrix-vector multiplication,  $O(m_{\bar{A}}p_{\bar{A}})$

$\mathbf{z} = \bar{\mathbf{A}}\mathbf{b}_i$

vector unpacking,  $O(n)$

for  $\tilde{i} = 1 : n$

$\hat{\mathbf{y}}_i(\tilde{i}) = \mathbf{z}(\mathbf{k}_{\bar{A}}(\tilde{i}))$

---

(a)  $\mathbf{X}_i\mathbf{b}_i = \bar{\mathbf{A}}\mathbf{b}_i$ , where  $\mathbf{X}_i$  is a singleton with packed storage  $(\bar{\mathbf{A}}, \mathbf{k}_{\bar{A}}, m_{\bar{A}}, p_{\bar{A}})$ . FLOP count for each computation step is given in gray text, and the algorithm involves  $O(n) + O(m_{\bar{A}}p_{\bar{A}})$  FLOP.

---

initialize a length- $n$  vector of zeros,  $O(n)$

$\hat{\mathbf{y}}_i = \text{zeros}(n)$

for  $k = 1 : p_{\bar{A}}$

discrete matrix-vector multiplication,  $O(n) + O(m_{\bar{A}s}p_{\bar{A}s})$

compute  $\mathbf{z} = \mathbf{A}_s\mathbf{b}_i((kp_{\bar{A}s} - p_{\bar{A}s} + 1) : kp_{\bar{A}s})$  using algorithm in Figure 8.6a

update  $\mathbf{z}$  and  $\hat{\mathbf{y}}_i$ ,  $O(sn)$

update  $\mathbf{z} = \text{diag}(\dot{\mathbf{A}}, k)\mathbf{z}$  using algorithm in Figure 8.3c

update  $\hat{\mathbf{y}}_i = \hat{\mathbf{y}}_i + \mathbf{z}$

---

(b)  $\mathbf{X}_i\mathbf{b}_i = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)\mathbf{b}_i$ , where  $\mathbf{X}_i$  is a tensor with  $s$  packed margins from  $(\bar{\mathbf{A}}_1, \mathbf{k}_{\bar{\mathbf{A}}_1}, m_{\bar{\mathbf{A}}_1}, p_{\bar{\mathbf{A}}_1})$  to  $(\bar{\mathbf{A}}_s, \mathbf{k}_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{A}}_s}, p_{\bar{\mathbf{A}}_s})$ , and  $\dot{\mathbf{A}} = \mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2 \tilde{\otimes} \cdots \tilde{\otimes} \mathbf{A}_{s-1}$ . FLOP count for each computation step is given in gray text, and the algorithm involves  $O(snp_{\dot{\mathbf{A}}}) + O(m_{\bar{\mathbf{A}}_s}p_{\dot{\mathbf{A}}}p_{\bar{\mathbf{A}}_s})$  FLOP.

---

**Figure 8.6:** Computation of  $\mathbf{X}_i\mathbf{b}_i$  when  $\mathbf{X}_i$  is a singleton or tensor.

thus an iterative algorithm as in Figure 8.6b applies, and it is  $O(\min\{p_{\bar{\mathbf{A}}_s}/s, n/m_{\bar{\mathbf{A}}_s}\})$  times faster than conventional computation.

With centring constraint, we need to compute  $\mathbf{X}_i\mathbf{Q}_i^\perp\mathbf{b}_i$ . By first obtaining  $\mathbf{b}_i^* = \mathbf{Q}_i^\perp\mathbf{b}_i$ ,  $\mathbf{X}_i\mathbf{b}_i^*$  can be computed with algorithms in Figure 8.6a or 8.6b.  $\mathbf{b}_i^*$  can be conveniently computed by  $\mathbf{b}_i^* = (\mathbf{Q}_i \mathbf{Q}_i^\perp) \begin{pmatrix} 0 \\ \mathbf{b}_i \end{pmatrix} = \mathbf{H}_1 \begin{pmatrix} 0 \\ \mathbf{b}_i \end{pmatrix}$ , where padding a zero to  $\mathbf{b}_i$  facilitates the use of Householder transformation.

## 8.6 Caching for computation of $\mathbf{X}'\mathbf{W}\mathbf{X}$

In §5.2.3 I have explained how a block algorithm or data caching is important for high-performance computing. It is possible to utilize data caching in the matrix weighting step of algorithms in Figures 8.4b, 8.5a and 8.5b to achieve higher performance when it comes to realistic computing on a modern computer with CPU cache. Common weight matrices include

- a diagonal matrix  $\mathbf{W} = \text{diag}(\mathbf{w})$  where only the main diagonal  $\mathbf{w}$  needs be stored;
- a symmetric tri-diagonal matrix, for example the precision matrix of an AR(1) process, like (2.1) or (6.1). Storage for such matrix only requires a vector  $\mathbf{w}$  of  $(2n - 1)$  elements by storing the main diagonal and a subdiagonal.

Let  $\mathbf{Y}$  be an  $n \times p_Y$  dense matrix, three types of matrix weighting will be met in these algorithms, namely

- matrix row-scaling that overwrites an existing matrix  $\mathbf{Y} = \text{diag}(\mathbf{w})\mathbf{Y}$ ;
- matrix row-scaling that writes to a new matrix  $\tilde{\mathbf{Y}} = \text{diag}(\mathbf{w})\mathbf{Y}$ ;

**Table 8.1:** Practical performance of three types of matrix weighting with and without caching. Experiment is taken on an Intel Core i5-2557M Sandy Bridge laptop (see Appendix A for hardware information) with a  $400000 \times 300$  matrix  $\mathbf{Y}$ , and execution time (in seconds) for 100 replications is given.

	$\mathbf{Y} = \text{diag}(\mathbf{w})\mathbf{Y}$	$\tilde{\mathbf{Y}} = \text{diag}(\mathbf{w})\mathbf{Y}$	$\mathbf{Y} = \mathbf{W}\mathbf{Y}$
no caching	33.6	41.4	61.7
caching	22.0	33.7	38.4

- symmetric tri-diagonal matrix weighting that overwrites an existing matrix  $\tilde{\mathbf{Y}} = \mathbf{W}\mathbf{Y}$ , where  $\mathbf{W}$  is a symmetric tri-diagonal matrix.

Table 8.1 is a practical example on how caching speeds up computation. In the following I will explain how caching is attained for these three types of matrix weighting.

Consider first  $\mathbf{Y} = \text{diag}(\mathbf{w})\mathbf{Y}$ . Conventional computation does column-wise scaling (see algorithm in Figure 8.7a). If  $n$  is small enough such that  $\mathbf{w}$  and  $\mathbf{Y}(:, j)$  fit in L1 CPU cache,  $\mathbf{w}$  only needs be fetched from RAM once, and it stays in cache until all columns of  $\mathbf{Y}$  are processed. The total number of reads from RAM is  $(np_Y + n)$ , including  $np_Y$  reads for  $\mathbf{Y}$  and  $n$  reads for  $\mathbf{w}$ . However, if  $n$  is big,  $\mathbf{w}$  has to be streamed from RAM for every column of  $\mathbf{Y}$ , as  $\mathbf{w}$  would be evicted from the cache after a column of  $\mathbf{Y}$  is processed. This results in  $2np_Y$  total number of reads from RAM. The amount of RAM access has a noticeable impact on performance, because the operation only involves  $np_Y$  floating-point multiplication, so it is memory-bound rather than computation-bound. To cache  $\mathbf{w}$  for arbitrary  $n$ , matrix row-scaling can be done tile-by-tile, as demonstrated by the algorithm in Figure 8.7b. In theory, tile size  $c$  should be chosen so that  $\mathbf{w}(h)$  occupies about half of the cache (because another half is reserved for a column of  $\mathbf{Y}$ ). Most computers today have 32KB L1 cache, which can hold 4096 double precision floating-point numbers. This means that  $c = 2000$  is near optimal. However, in practice it seems that performance is not very sensitive to the exact choice of  $c$ ; any values between 200 and 4000 give about the same performance. But if for example,  $c = 100$ , the execution time becomes 26.4 seconds.

For  $\tilde{\mathbf{Y}} = \text{diag}(\mathbf{w})\mathbf{Y}$ , the same kind of tiling can be done, except that the theoretically optimal tile size is chosen so that  $\mathbf{w}(h)$  occupies 1/3 of cache (because a column of  $\tilde{\mathbf{Y}}$  and  $\mathbf{Y}$  both require 1/3 of cache, too), giving  $c = 1350$ . For this operation, performance is more sensitive to  $c$ . If for example,  $c = 800$  or  $c = 1600$ , execution time becomes about 36.6 seconds.

Now consider  $\mathbf{Y} = \mathbf{W}\mathbf{Y}$ , or

$$\begin{pmatrix} \mathbf{w}(1) & \mathbf{w}(2) & & & \\ \mathbf{w}(2) & \mathbf{w}(3) & \mathbf{w}(4) & & \\ & \ddots & \ddots & \ddots & \\ & & \mathbf{w}(2n-4) & \mathbf{w}(2n-3) & \mathbf{w}(2n-2) \\ & & & \mathbf{w}(2n-2) & \mathbf{w}(2n-1) \end{pmatrix} \begin{pmatrix} \cdots & \mathbf{Y}(:, 1) & \cdots \\ \cdots & \mathbf{Y}(:, 2) & \cdots \\ & \vdots & \\ \cdots & \mathbf{Y}(:, n-1) & \cdots \\ \cdots & \mathbf{Y}(:, n) & \cdots \end{pmatrix},$$

with a packed storage for  $\mathbf{W}$ . Conventional computation does column-wise matrix-vector multiplication (see algorithm in Figure 8.7c), but if  $n$  is too big for  $\mathbf{w}$  and  $\mathbf{Y}(:, j)$  to fit in cache,  $\mathbf{w}$  has to be streamed from RAM for every column of  $\mathbf{Y}$ . We can cache  $\mathbf{w}$  by row tiling, but a working vector  $\mathbf{u}$  is needed as a buffer to store a row of  $\mathbf{Y}$ . Algorithm in Figure 8.7d implements this. Caching has a even more pronounced effect on this operation than matrix row-scaling, because it reduces total amount of RAM access from  $3np_Y$  to  $(np_Y + 2n - 1)$ . As  $(2c - 1)$  elements of  $\mathbf{w}$  and  $c$  elements of  $\mathbf{Y}$  would be accessed in the innermost  $i$ -loop,  $c$  should be chosen so that  $3c$  double precision floating-point numbers can fit in cache. With a 32KB L1 cache, this is  $c = 1350$ . But once again, performance is insensitive to  $c$  in practice, as long as it is not too small (say 100) or too big (say 4000).

for $j = 1 : p_Y$ $\mathbf{Y}(, j) = \mathbf{w} \circ \mathbf{Y}(, j)$	
(a) uncached matrix row-scaling $\mathbf{Y} = \text{diag}(\mathbf{w})\mathbf{Y}$	
$i = 1$ while $(i \leq n)$ $\tilde{c} = \min\{c, n - i + 1\}$ row-scaling of $\mathbf{Y}(i : (i + \tilde{c} - 1), )$ $\mathbf{Y}(i : (i + \tilde{c} - 1), ) = \text{diag}(\mathbf{w}(i : (i + \tilde{c} - 1)))\mathbf{Y}(i : (i + \tilde{c} - 1), )$ $i = i + \tilde{c}$	
(b) cached matrix row-scaling $\mathbf{Y} = \text{diag}(\mathbf{w})\mathbf{Y}$	
for $j = 1 : p_Y$ update $\mathbf{Y}(1, j)$ $r = \mathbf{Y}(1, j)$ $\mathbf{Y}(1, j) = \mathbf{w}(1) \cdot r + \mathbf{w}(2) \cdot \mathbf{Y}(2, j)$ update $\mathbf{Y}(2 : (n - 1), j)$ for $i = 2 : (n - 1)$ $s = \mathbf{w}(2i - 2) \cdot r + \mathbf{w}(2i - 1) \cdot \mathbf{Y}(i, j) + \mathbf{w}(2i) \cdot \mathbf{Y}(i + 1, j)$ $r = \mathbf{Y}(i, j)$ $\mathbf{Y}(i, j) = s$ update $\mathbf{Y}(n, j)$ $\mathbf{Y}(n, j) = \mathbf{w}(2n - 2) \cdot r + \mathbf{w}(2n - 1) \cdot \mathbf{Y}(n, j)$	
(c) uncached tri-diagonal weighting $\mathbf{Y} = \mathbf{W}\mathbf{Y}$	
initialize a length- $p_Y$ buffer vector $\mathbf{u}$ $\mathbf{u} = \text{zeros}(p_Y)$ update $\mathbf{Y}(1, )$ for $j = 1 : p_Y$ $\mathbf{u}(j) = \mathbf{Y}(1, j)$ $\mathbf{Y}(1, j) = \mathbf{w}(1) \cdot \mathbf{u}(j) + \mathbf{w}(2) \cdot \mathbf{Y}(2, j)$ update $\mathbf{Y}(2 : (n - 1), j)$ , one tile per time $i = 2$ while $(i \leq (n - 1))$ $\tilde{c} = \min\{c, n - i\}$ for $j = 1 : p_Y$ $r = \mathbf{u}(j)$ for $\tilde{i} = i : (i + \tilde{c} - 1)$ $s = \mathbf{w}(2\tilde{i} - 2) \cdot r + \mathbf{w}(2\tilde{i} - 1) \cdot \mathbf{Y}(\tilde{i}, j) + \mathbf{w}(2\tilde{i}) \cdot \mathbf{Y}(\tilde{i} + 1, j)$ $r = \mathbf{Y}(\tilde{i}, j)$ $\mathbf{Y}(\tilde{i}, j) = s$ $\mathbf{u}(j) = r$ $i = i + \tilde{c}$ update $\mathbf{Y}(n, )$ for $j = 1 : p_Y$ $\mathbf{Y}(n, j) = \mathbf{w}(2n - 2) \cdot \mathbf{u}(j) + \mathbf{w}(2n - 1) \cdot \mathbf{Y}(n, j)$	
(d) cached tri-diagonal weighting $\mathbf{Y} = \mathbf{W}\mathbf{Y}$	

**Figure 8.7:** Algorithms for matrix weighting with and without caching.

## 8.7 Experiment on daily logBS model

To demonstrate that algorithms elaborated in previous sections yield significant speedup to pseudo QR reduction and  $\hat{\mathbf{y}}$  computation in practice, consider their application to the test model (model 7.1) with 1342159 data and 22571 coefficients. A quick test with an arbitrary length-1342159 vector  $\mathbf{y}$  and an arbitrary length-22571 vector  $\hat{\boldsymbol{\beta}}_{\lambda}$  shows that computations of  $\mathbf{X}'\mathbf{W}\mathbf{y}$  and  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}_{\lambda}$  respectively take 1.5 and 2.5 seconds in serial computing. By contrast, computing  $\hat{\mathbf{y}}$  with “online” algorithm and the fast model matrix formation method described in §8.1 still take 3 minutes using 16 threads.

Computation of  $\mathbf{X}'\mathbf{W}\mathbf{X}$  takes 11 hours with single-threaded OpenBLAS if using the “online” pseudo QR reduction algorithm. With 16-threaded OpenBLAS it takes 50 minutes. Discrete computation of  $\mathbf{X}'\mathbf{W}\mathbf{X}$  takes 5 hours 40 minutes in serial computing, and 34 minutes in 16-threaded parallel computing. Here, multi-threaded parallel computing for  $\mathbf{X}'\mathbf{W}\mathbf{X}$  is achieved by parallelizing all steps of the algorithms via OpenMP.

- The matrix row-aggregation in Figure 8.4a can process different columns in parallel;
- With the caching algorithm in Figure 8.7b, the matrix row-scaling in Figure 8.4b can process different row tiles in parallel;
- The column unpacking of  $\mathbf{B}$  in Figure 8.5a and 8.5b can process different columns in parallel;
- The matrix weighting in Figure 8.5a and 8.5b is also parallelizable. For a diagonal  $\mathbf{W}$ , this is just a parallel matrix row-scaling. For a symmetric tri-diagonal  $\mathbf{W}$ , the caching algorithm in Figure 8.7d implies that there is dependency between different row tiles through the buffer vector. However, we can instead split  $\mathbf{Y}$  into a few column chunks, and weight different chunks in parallel.

All these computations are memory-bound, so parallel scaling is up to how well different CPUs in a machine can do parallel read from RAM or parallel write to RAM. Different machines may have different front-side bus (FSB) bandwidth, thus the scaling factor for these operations are likely to vary from machine to machine. On an Intel Xeon E5-2650 v2 workstation, parallel scaling is not too bad, with 5.5 times speedup at 8 threads and 10 times speedup at 16 threads.

## 8.8 Summary

This Chapter has developed algorithms for computing matrix cross-product  $\mathbf{X}'\mathbf{W}\mathbf{X}$  when the sub-matrices of  $\mathbf{X}$  corresponding to individual GAM terms in a linear predictor  $\boldsymbol{\eta} = \mathbf{X}_0\boldsymbol{\beta}_0 + \sum_{i=1}^l \mathbf{X}_i\boldsymbol{\beta}_i$  are stored in a packed storage format from covariate discretization. Computation of each  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_j$  is  $O(p)$  times faster than the “online” pseudo QR reduction algorithm, where  $p$  is the number of columns in  $\mathbf{X}_i$  if  $\mathbf{X}_i$  is a singleton, and that of the last margin of  $\mathbf{X}_i$  if  $\mathbf{X}_i$  is a tensor.

This idea of covariate discretization was first introduced in Lang et al. (2014) to obtain efficient storage and computation for large datasets. However, in that paper they employ smooths of one covariate and only require terms of the form  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_i$ , but not  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_j$ . So methods in this Chapter have extended the idea by

- discretizing covariates in multivariate splines (say tensor product splines) and exploiting the resulting packed storage for computation;
- developing algorithms for computing any  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_j$ .



Regarding the first point, method in this Chapter discretizes covariates marginally rather than jointly, so that spline functions can be evaluated at their unique combinations. The discretization result is also stored marginally rather than jointly, which substantially reduces the storage complexity. This is a major advantage over Helwig and Ma (2016). For example, consider a two margin tensor product spline. To store  $s$  unique covariate values of the first margin and  $t$  unique covariate values of the second margin, the marginal approach only requires  $(s + t)$  storage costs, but the joint approach requires  $st$  storage costs. As a result, the joint approach has more difficulty in balancing storage costs and precision of discretization.

A hard question is how to interpret “Computation of each  $\mathbf{X}_i' \mathbf{W} \mathbf{X}_j$  is  $O(p)$  times faster than the “online” pseudo QR reduction algorithm”. Indeed, this does not seem to agree with the benchmarking result in §8.7, where the new discrete computation method is only two times as fast as the old “online” algorithm. The problem is, that  $O(p)$  factor is derived from comparing FLOP count of the two computation methods, but in reality FLOP count is not everything. The “online” pseudo QR reduction, while computationally more expensive, can attain high performance computing with the aid of an optimized BLAS. The new discrete method, while computationally cheaper, is not rich in level-3 BLAS operations. On the contrary, it is rich in vector-oriented operations like unpacking, aggregation, reweighting, etc. These operations are like level-2 BLAS operations; they are memory-bound rather than CPU-bound. In §8.6, some attempts were made to improve the performance of this new computation method, however, just like a matrix-vector multiplication can hardly be boosted by a block algorithm or data caching (see a relevant discussion on this in §5.2.3), the practical gains is very limited.

Is there any way, to overcome this drawback of the discrete computation method? Yes, there is. It is **bamboos**, a new computational engine for discrete  $\mathbf{X}' \mathbf{W} \mathbf{X}$ . Welcome to the next Chapter, the most exciting Chapter of this thesis.

## Chapter 9

# Bamboos: boosting discrete pseudo QR reduction by another magnitude

In this Chapter, I will introduce **bamboos** (**bam** booster), a new computational engine for matrix cross-product  $\mathbf{X}'\mathbf{W}\mathbf{X}$  that is truly high performance computing, able to boost the discrete computation method developed in the previous Chapter by another magnitude. This Chapter is laid out as follows.

- §9.1 explains the basic idea of **bamboos**;
- §9.2 makes FLOP count comparison between **bamboos** and the old discrete algorithms;
- §9.3 explains how **bamboos** actually achieves high performance;
- §9.4 benchmarks **bamboos** and the old discrete algorithms by fitting the test daily logBS model, model 7.1;
- §9.5 summarizes key features and advantages of **bamboos**.

### 9.1 A basic description of bamboos algorithms

Recall that for a linear predictor  $\boldsymbol{\eta} = \mathbf{X}_0\boldsymbol{\beta}_0 + \sum_{i=1}^l \mathbf{f}_i = \mathbf{X}_0\boldsymbol{\beta}_0 + \sum_{i=1}^l \mathbf{X}_i\boldsymbol{\beta}_i$  with a design matrix  $\mathbf{X}_0$  for parametric model terms and design matrices  $\mathbf{X}_i$ 's for splines  $\mathbf{f}_i$ 's, algorithms in §8.4 can be interpreted as such:

- For any  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_0$ , treat  $\mathbf{X}_0$  as  $\mathbf{Y}$ , then compute  $\mathbf{X}_i'\mathbf{W}\mathbf{Y}$  using algorithms in Figures 8.4a and 8.4b;
- for any  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_j$  ( $i \geq j > 0$ ), unpack  $\mathbf{X}_j$  from its packed storage into a full design matrix  $\mathbf{Y}$ , then compute  $\mathbf{X}_i'\mathbf{W}\mathbf{Y}$  using algorithms in Figures 8.4a and 8.4b.

So computation of  $\mathbf{X}_i'\mathbf{W}\mathbf{Y}$  is elementary. By contrast, in **bamboos**, such computation is only useful for computing  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_0$ . For other blocks, computation of  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_j$  between two singletons plays a pivotal role, as  $\mathbf{X}_i'\mathbf{W}\mathbf{X}_j$  involving tensors would be cast into its iterative application. As it will be shown, there exists algorithms for computing a “singleton  $\times$  singleton” block without matrix unpacking, hence **bamboos** algorithms require no matrix unpacking at all! In the rest of this section,

- §9.1.1 motivates algorithms for “singleton  $\times$  singleton” blocks;

- §9.1.2 extends these algorithms for “singleton  $\times$  Tensor”, “Tensor  $\times$  singleton” and “Tensor  $\times$  Tensor” blocks;
- §9.1.3 is a consolidation, summarizing above algorithms with a unified representation.

In the forthcoming demonstration, I will adopt the same notations used in Chapter 8. For singletons  $\mathbf{X}_i$  and  $\mathbf{X}_j$ , denote that they respectively have packed storage  $(\bar{\mathbf{A}}, \mathbf{k}_{\bar{\mathbf{A}}}, m_{\bar{\mathbf{A}}}, p_{\bar{\mathbf{A}}})$  and  $(\bar{\mathbf{B}}, \mathbf{k}_{\bar{\mathbf{B}}}, m_{\bar{\mathbf{B}}}, p_{\bar{\mathbf{B}}})$ . For tensors  $\mathbf{X}_i$  and  $\mathbf{X}_j$ , denote that they respectively have  $s$  packed margins from  $(\bar{\mathbf{A}}_1, \mathbf{k}_{\bar{\mathbf{A}}_1}, m_{\bar{\mathbf{A}}_1}, p_{\bar{\mathbf{A}}_1})$  to  $(\bar{\mathbf{A}}_s, \mathbf{k}_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{A}}_s}, p_{\bar{\mathbf{A}}_s})$ , and  $t$  packed margins from  $(\bar{\mathbf{B}}_1, \mathbf{k}_{\bar{\mathbf{B}}_1}, m_{\bar{\mathbf{B}}_1}, p_{\bar{\mathbf{B}}_1})$  to  $(\bar{\mathbf{B}}_t, \mathbf{k}_{\bar{\mathbf{B}}_t}, m_{\bar{\mathbf{B}}_t}, p_{\bar{\mathbf{B}}_t})$ . I also define  $\tilde{\mathbf{A}} = \mathbf{A}_1 \tilde{\otimes} \mathbf{A}_2 \tilde{\otimes} \cdots \mathbf{A}_{s-1}$  with  $p_{\tilde{\mathbf{A}}} = \prod_{i=1}^{s-1} p_{\bar{\mathbf{A}}_i}$ , and  $\tilde{\mathbf{B}} = \mathbf{B}_1 \tilde{\otimes} \mathbf{B}_2 \tilde{\otimes} \cdots \mathbf{B}_{t-1}$  with  $p_{\tilde{\mathbf{B}}} = \prod_{i=1}^{t-1} p_{\bar{\mathbf{B}}_i}$ , so that any tensor can be expressed with two margins.

### 9.1.1 $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$ between singletons

Consider first a diagonal weight matrix  $\mathbf{W} = \text{diag}(\mathbf{w})$ , then  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_i = \mathbf{A}' \mathbf{W} \mathbf{A}$  and  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j = \mathbf{A}' \mathbf{W} \mathbf{B}$  can be computed using algorithms in Figure 9.1a. Note that unlike the discrete algorithms in the last Chapter where weights are absorbed into the design matrix on the right and processed together with a matrix, here, weights are processed independently and are aggregated into a contingency table / matrix  $\bar{\mathbf{W}}$ . Below is a mathematical proof of the algorithm for  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j = \mathbf{A}' \mathbf{W} \mathbf{B}$ . Partition  $\mathbf{A}$  and  $\mathbf{B}$  by rows, it can then be recognized that  $\mathbf{A}' \mathbf{W} \mathbf{B}$  is a quadratic form:

$$\mathbf{A}' \mathbf{W} \mathbf{B} = \begin{bmatrix} \mathbf{A}'(1, ) & \mathbf{A}'(2, ) & \cdots & \mathbf{A}'(n, ) \end{bmatrix} \begin{bmatrix} w_1 & & & \\ & w_2 & & \\ & & \ddots & \\ & & & w_n \end{bmatrix} \begin{bmatrix} \mathbf{B}(1, ) \\ \mathbf{B}(2, ) \\ \vdots \\ \mathbf{B}(n, ) \end{bmatrix} = \sum_{i=1}^n w_i \mathbf{A}'(i, ) \mathbf{B}(i, ).$$

Now I am going to tweak the summation to sum over unique rows of  $\mathbf{A}$  and  $\mathbf{B}$ , i.e., rows of packed design matrices  $\bar{\mathbf{A}}$  and  $\bar{\mathbf{B}}$ . Essentially it is grouping  $w_i$  by  $(\mathbf{k}_{\bar{\mathbf{A}}}(i), \mathbf{k}_{\bar{\mathbf{B}}}(i))$  pairs, giving:

$$\sum_{u=1}^{p_{\bar{\mathbf{A}}}} \sum_{v=1}^{p_{\bar{\mathbf{B}}}} \left( \sum_{\substack{\mathbf{k}_{\bar{\mathbf{A}}}(i)=u \\ \mathbf{k}_{\bar{\mathbf{B}}}(i)=v}} w_i \right) \bar{\mathbf{A}}'(u, ) \bar{\mathbf{B}}(v, ) := \sum_{u=1}^{p_{\bar{\mathbf{A}}}} \sum_{v=1}^{p_{\bar{\mathbf{B}}}} \bar{w}_{uv} \bar{\mathbf{A}}'(u, ) \bar{\mathbf{B}}(v, ) = \bar{\mathbf{A}}' \bar{\mathbf{W}} \bar{\mathbf{B}},$$

where  $\bar{\mathbf{W}}(u, v) = w_{uv}$  is the group sum, which is what the in the algorithm computes.

The advantage of these algorithms is most obvious for  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_i$ , whose leading order FLOP count is  $O(n) + O(m_{\bar{\mathbf{A}}} p_{\bar{\mathbf{A}}}^2)$ . By contrast, direct computation that first unpacks  $\bar{\mathbf{A}}$  for  $\mathbf{A}$ , then performs row-scaling  $\tilde{\mathbf{A}} = \mathbf{W} \mathbf{A}$  and finally computes matrix cross-product  $\mathbf{A}' \tilde{\mathbf{A}}$  involves  $O(np_{\bar{\mathbf{A}}}) + O(np_{\bar{\mathbf{A}}}) + O(np_{\bar{\mathbf{A}}}^2)$  FLOP. So the algorithm is very efficient when  $n \gg m_{\bar{\mathbf{A}}}$ .

For  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$ , the matrix product and cross-product can be computed in two different orders, and the least expensive one should be chosen in practice. But the advantage of the algorithm is not immediately clear, as while  $n \gg m_{\bar{\mathbf{A}}}$  and  $n \gg m_{\bar{\mathbf{B}}}$  may hold, it is not guaranteed that  $n \gg m_{\bar{\mathbf{A}}} m_{\bar{\mathbf{B}}}$ . For example, consider  $n = 10^5$  with  $m_{\bar{\mathbf{A}}} = m_{\bar{\mathbf{B}}} = 10^3$ . So the  $O(m_{\bar{\mathbf{A}}} m_{\bar{\mathbf{B}}} p_{\bar{\mathbf{B}}})$  or  $O(m_{\bar{\mathbf{A}}} m_{\bar{\mathbf{B}}} p_{\bar{\mathbf{A}}})$  complexity for the matrix product can be very expensive, as  $m_{\bar{\mathbf{A}}} m_{\bar{\mathbf{B}}}$  approaches or even goes beyond  $n$ . However, since  $\bar{\mathbf{W}}$  is obtained from  $n$  data, it at most has  $n$  non-zero entries. As  $m_{\bar{\mathbf{A}}} m_{\bar{\mathbf{B}}}$  becomes increasingly large,  $\bar{\mathbf{W}}$  becomes increasingly sparse, so I can avoid forming  $\bar{\mathbf{W}}$  as a dense matrix and do sparse matrix multiplication instead, which only involves  $O(K p_{\bar{\mathbf{B}}})$  or  $O(K p_{\bar{\mathbf{A}}})$  FLOP, where  $K$  is the number of non-zero entries in  $\bar{\mathbf{W}}$  (sparse  $\bar{\mathbf{W}}$  will be covered in details in §9.3). Obviously,  $K < \min\{m_{\bar{\mathbf{A}}} m_{\bar{\mathbf{B}}}, n\}$ , so its computational complexity is bounded. In the worst case, the algorithm has a leading order FLOP count of  $O(n) + O(np_{\bar{\mathbf{B}}}) + O(m_{\bar{\mathbf{A}}} p_{\bar{\mathbf{A}}} p_{\bar{\mathbf{B}}})$  or  $O(n) + O(np_{\bar{\mathbf{A}}}) + O(m_{\bar{\mathbf{B}}} p_{\bar{\mathbf{A}}} p_{\bar{\mathbf{B}}})$ , yet this is still a lot more efficient than direct computation that first unpacks  $\bar{\mathbf{A}}$  and  $\bar{\mathbf{B}}$  to  $\mathbf{A}$  and  $\mathbf{B}$  (with  $O(np_{\bar{\mathbf{A}}}) + O(np_{\bar{\mathbf{B}}})$  FLOP), then performs row-scaling  $\tilde{\mathbf{A}} = \mathbf{W} \mathbf{A}$  or  $\tilde{\mathbf{B}} = \mathbf{W} \mathbf{B}$  (with  $O(np_{\bar{\mathbf{A}}})$  or  $O(np_{\bar{\mathbf{B}}})$  FLOP).

initialize a length- $m_{\bar{A}}$ vector of zeros $\bar{\mathbf{w}} = \text{zeros}(m_{\bar{A}})$ one-way contingency table, $O(n)$ for $i = 1 : n$ $\bar{\mathbf{w}}(\mathbf{k}_{\bar{A}}(i)) += \mathbf{w}(i)$ matrix row-scaling, $O(m_{\bar{A}}p_{\bar{A}})$ $\tilde{\mathbf{A}} = \text{diag}(\bar{\mathbf{w}})\mathbf{\bar{A}}$ matrix cross-product, $O(m_{\bar{A}}p_{\bar{A}}^2)$ $\mathbf{X}'_i \mathbf{W} \mathbf{X}_i = \mathbf{\bar{A}}' \tilde{\mathbf{A}}$	initialize a $m_{\bar{A}} \times m_{\bar{B}}$ matrix of zeros $\bar{\mathbf{W}} = \text{zeros}(m_{\bar{A}}, m_{\bar{B}})$ two-way contingency table, $O(n)$ for $i = 1 : n$ $\bar{\mathbf{W}}(\mathbf{k}_{\bar{A}}(i), \mathbf{k}_{\bar{B}}(i)) += \mathbf{w}(i)$ matrix product, $O(m_{\bar{A}}m_{\bar{B}}p_{\bar{B}})$ or $O(m_{\bar{A}}m_{\bar{B}}p_{\bar{A}})$ $\tilde{\mathbf{B}} = \bar{\mathbf{W}}\bar{\mathbf{B}}$ or $\tilde{\mathbf{A}} = \mathbf{\bar{A}}'\bar{\mathbf{W}}$ cross-product, $O(m_{\bar{A}}p_{\bar{A}}p_{\bar{B}})$ or $O(m_{\bar{B}}p_{\bar{A}}p_{\bar{B}})$ $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j = \mathbf{\bar{A}}' \tilde{\mathbf{B}}$ or $\tilde{\mathbf{A}} \tilde{\mathbf{B}}$
--	--

(a) Computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_i = \mathbf{A}' \mathbf{W} \mathbf{A}$  and  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j = \mathbf{A}' \mathbf{W} \mathbf{B}$ , where  $\mathbf{W} = \text{diag}(\mathbf{w})$  is a diagonal matrix. Note that if matrix row-scaling is performed by  $\tilde{\mathbf{A}} = \text{diag}(\sqrt{\mathbf{w}})\mathbf{\bar{A}}$ , the matrix cross-product  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_i = \tilde{\mathbf{A}}' \tilde{\mathbf{A}}$  will involve 50% less FLOP count. However, This idea does not work when there are negative weights in  $\mathbf{w}$ . This can happen in two cases: 1) the working weights in P-IRLS (see algorithm in Figure 4.12) can contain negative values; 2) the weights can be negative after being reweighted (see algorithms in the next section).

---

initialize a  $m_{\bar{A}} \times m_{\bar{B}}$  matrix of zeros  
 $\bar{\mathbf{W}} = \text{zeros}(m_{\bar{A}}, m_{\bar{B}})$   
 two-way contingency table,  $O(n)$   
 for  $i = 1 : n$   
      $\bar{\mathbf{W}}(\mathbf{k}_{\bar{A}}(i), \mathbf{k}_{\bar{B}}(i)) += \mathbf{w}_0(i)$   
 for  $i = 1 : (n - 1)$   
      $\bar{\mathbf{W}}(\mathbf{k}_{\bar{A}}(i), \mathbf{k}_{\bar{B}}(i + 1)) += \mathbf{w}_1(i)$   
 for  $i = 1 : (n - 1)$   
      $\bar{\mathbf{W}}(\mathbf{k}_{\bar{A}}(i + 1), \mathbf{k}_{\bar{B}}(i)) += \mathbf{w}_{-1}(i)$   
 matrix product,  $O(m_{\bar{A}}m_{\bar{B}}p_{\bar{B}})$  or  $O(m_{\bar{A}}m_{\bar{B}}p_{\bar{A}})$   
 $\tilde{\mathbf{B}} = \bar{\mathbf{W}}\bar{\mathbf{B}}$  or  $\tilde{\mathbf{A}} = \mathbf{\bar{A}}'\bar{\mathbf{W}}$   
 matrix cross-product,  $O(m_{\bar{A}}p_{\bar{A}}p_{\bar{B}})$  or  $O(m_{\bar{B}}p_{\bar{A}}p_{\bar{B}})$   
 $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j = \mathbf{\bar{A}}' \tilde{\mathbf{B}}$  or  $\tilde{\mathbf{A}} \tilde{\mathbf{B}}$

---

(b) Computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j = \mathbf{A}' \mathbf{W} \mathbf{B}$ , where  $\mathbf{W}$  is a tri-diagonal matrix (either symmetric or asymmetric) whose main diagonal, upper subdiagonal and lower subdiagonal are stored as  $\mathbf{w}_0$ ,  $\mathbf{w}_1$  and  $\mathbf{w}_{-1}$  respectively. See main text on why computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_i$  is not provided.

**Figure 9.1:** Computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$  when both  $\mathbf{X}_i$  and  $\mathbf{X}_j$  are singletons, respectively with packed storage  $(\bar{\mathbf{A}}, \mathbf{k}_{\bar{A}}, m_{\bar{A}}, p_{\bar{A}})$  and  $(\bar{\mathbf{B}}, \mathbf{k}_{\bar{B}}, m_{\bar{B}}, p_{\bar{B}})$ . FLOP count for each computation step is given in gray text. Note that the result for matrix product is given assuming that  $\bar{\mathbf{W}}$  is (constructed as) dense, therefore the FLOP count is proportional to its size  $m_{\bar{A}}m_{\bar{B}}$ . But  $\bar{\mathbf{W}}$  can be sparse with  $K < \min\{m_{\bar{A}}m_{\bar{B}}, n\}$  (or even  $K \ll \min\{m_{\bar{A}}m_{\bar{B}}, n\}$ ) non-zero elements. Constructing  $\bar{\mathbf{W}}$  as sparse and exploiting sparsity in matrix-matrix multiplication in this situation can substantially reduce the  $O(m_{\bar{A}}m_{\bar{B}}p_{\bar{B}})$  or  $O(m_{\bar{A}}m_{\bar{B}}p_{\bar{A}})$  complexity to  $O(Kp_{\bar{B}})$  or  $O(Kp_{\bar{A}})$ . Details on sparse construction of  $\bar{\mathbf{W}}$  will be given later in §9.3.

and finally computes the matrix cross-product  $\tilde{\mathbf{A}}' \tilde{\mathbf{B}}$  or  $\mathbf{\bar{A}}' \tilde{\mathbf{B}}$  (with  $O(np_{\bar{A}}p_{\bar{B}})$  FLOP) when  $n \gg m_{\bar{A}}$  and  $n \gg m_{\bar{B}}$ .

When  $\mathbf{W}$  is a tri-diagonal matrix like (2.1) or (6.1), it can be decomposed into three matrices:

$$\begin{pmatrix} \bullet & & & & \\ \bullet & \bullet & & & \\ & \bullet & \bullet & & \\ & & \bullet & \bullet & \\ & & & \bullet & \bullet \end{pmatrix} = \begin{pmatrix} \bullet & & & & \\ & \bullet & & & \\ & & \bullet & & \\ & & & \bullet & \\ & & & & \bullet \end{pmatrix} + \begin{pmatrix} 0 & \bullet & & & \\ 0 & & \bullet & & \\ 0 & & & \bullet & \\ 0 & & & & \bullet \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ \bullet & & & & \\ & \bullet & & & \\ & & \bullet & & \\ & & & \bullet & \\ & & & & \bullet \end{pmatrix}.$$

Denote the first diagonal matrix by  $\mathbf{W}_0$ , and the diagonal submatrices in the second and third matrices by  $\mathbf{W}_1$  and  $\mathbf{W}_{-1}$ , respectively. Then for example,  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$  decomposes into

$$\mathbf{X}'_i \mathbf{W}_0 \mathbf{X}_j + \mathbf{X}_i(1 : (n - 1), )' \mathbf{W}_1 \mathbf{X}_j(2 : n, ) + \mathbf{X}_i(2 : n, )' \mathbf{W}_{-1} \mathbf{X}_j(1 : (n - 1), ),$$

where each part can be computed using algorithms in Figure 9.1a. See Figure 9.1b for an example implementation (and note that the three s can be merged for conciseness). Computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_i =$

---

for  $k = 1 : p_{\dot{A}}$

compute  $\mathbf{a} = \dot{\mathbf{A}}(\cdot, k)$  using algorithm in Figure 8.3c

$\mathbf{W}_k = \text{diag}(\mathbf{a})\mathbf{W}$

compute  $\mathbf{Z} = \mathbf{A}'_s \mathbf{W}_k \mathbf{B}$  using algorithms in Figure 9.1a or 9.1b

$(\mathbf{X}'_i \mathbf{W} \mathbf{X}_j)((kp_{\bar{A}_s} - p_{\bar{A}_s} + 1) : kp_{\bar{A}_s}, \cdot) = \mathbf{Z}$

---

(a) Computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} \mathbf{B}$ , where  $\mathbf{X}_j$  is a singleton with packed storage  $(\bar{\mathbf{B}}, \mathbf{k}_{\bar{\mathbf{B}}}, m_{\bar{\mathbf{B}}}, p_{\bar{\mathbf{B}}})$ .

---

for  $k = 1 : p_{\dot{A}}$

compute  $\mathbf{a} = \dot{\mathbf{A}}(\cdot, k)$  using algorithm in Figure 8.3c

$\mathbf{W}_k = \text{diag}(\mathbf{a})\mathbf{W}$

compute  $\mathbf{Z} = (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)' \mathbf{W}'_k \mathbf{A}_s$  using algorithm in Figure 9.2a

$(\mathbf{X}'_i \mathbf{W} \mathbf{X}_j)((kp_{\bar{A}_s} - p_{\bar{A}_s} + 1) : kp_{\bar{A}_s}, \cdot) = \mathbf{Z}'$

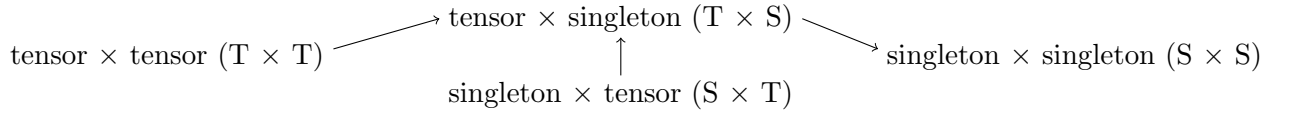
---

(b) Computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)$ , where  $\mathbf{X}_j$  is a tensor with  $t$  packed margins from  $(\bar{\mathbf{B}}_1, \mathbf{k}_{\bar{\mathbf{B}}_1}, m_{\bar{\mathbf{B}}_1}, p_{\bar{\mathbf{B}}_1})$  to  $(\bar{\mathbf{B}}_t, \mathbf{k}_{\bar{\mathbf{B}}_t}, m_{\bar{\mathbf{B}}_t}, p_{\bar{\mathbf{B}}_t})$ .

---

**Figure 9.2:** Computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$  when  $\mathbf{X}_i$  is a tensor with  $s$  packed margins from  $(\bar{\mathbf{A}}_1, \mathbf{k}_{\bar{\mathbf{A}}_1}, m_{\bar{\mathbf{A}}_1}, p_{\bar{\mathbf{A}}_1})$  to  $(\bar{\mathbf{A}}_s, \mathbf{k}_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{A}}_s}, p_{\bar{\mathbf{A}}_s})$ , and  $\mathbf{X}_j$  is a singleton or a tensor.

---



**Figure 9.3:** Illustrating the computational idea of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$  with a dependency graph. Matrix cross-product between two singletons is the elementary computation, and other cases are its iterative application.

---

$\mathbf{A}' \mathbf{W} \mathbf{A}$  can simply use the same algorithm. The seemingly more efficient computations by first doing a Cholesky factorization  $\bar{\mathbf{W}} = \mathbf{U}' \mathbf{U}$ , then computing a triangular matrix multiplication  $\tilde{\mathbf{A}} = \mathbf{U} \bar{\mathbf{A}}$ , and finally taking a matrix cross-product  $\mathbf{X}_i \mathbf{W} \mathbf{X}_i = \tilde{\mathbf{A}}' \mathbf{W} \tilde{\mathbf{A}}$  does not generally apply, because  $\mathbf{W}$  may be asymmetric as will be seen in the next section.

### 9.1.2 $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$ with tensors

Assume that  $\mathbf{X}_i$  is a tensor, then by applying (1.1) to  $\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s$ , there is

$$(\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} \mathbf{X}_j = \begin{bmatrix} \mathbf{A}'_s \text{diag}(\dot{\mathbf{A}}(\cdot, 1)) \\ \mathbf{A}'_s \text{diag}(\dot{\mathbf{A}}(\cdot, 2)) \\ \vdots \\ \mathbf{A}'_s \text{diag}(\dot{\mathbf{A}}(\cdot, p_{\dot{\mathbf{A}}})) \end{bmatrix} \mathbf{W} \mathbf{X}_j = \begin{bmatrix} \mathbf{A}'_s (\text{diag}(\dot{\mathbf{A}}(\cdot, 1)) \mathbf{W} \mathbf{X}_j) \\ \mathbf{A}'_s (\text{diag}(\dot{\mathbf{A}}(\cdot, 2)) \mathbf{W} \mathbf{X}_j) \\ \vdots \\ \mathbf{A}'_s (\text{diag}(\dot{\mathbf{A}}(\cdot, p_{\dot{\mathbf{A}}})) \mathbf{W} \mathbf{X}_j) \end{bmatrix}.$$

That is,  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$  can be obtained by computing a matrix cross-product  $\mathbf{A}'_s \mathbf{W}_k \mathbf{X}_j$  block by block, where the  $k^{\text{th}}$  block has a working weight matrix  $\mathbf{W}_k = \text{diag}(\dot{\mathbf{A}}(\cdot, k)) \mathbf{W}$ . Be aware that when  $\mathbf{W}$  is symmetric tri-diagonal,  $\mathbf{W}_k$  is asymmetric. Now,

1. if  $\mathbf{X}_j$  is a singleton, computation of this block can use algorithms in Figure 9.1a or 9.1b, giving an iterative algorithm in Figure 9.2a;
2. if  $\mathbf{X}_j$  is also a tensor, I can first compute  $\mathbf{X}'_j \mathbf{W}'_k \mathbf{A}_s$  using the previous algorithm, then transpose the result. This gives an iterative algorithm in Figure 9.2b.

In summary, such computational idea is illustrated by the dependency graph in Figure 9.3.

---

```

for  $h = 1 : p_{\dot{B}}$ 
  compute  $\mathbf{b} = \mathbf{B}(\cdot, h)$  using algorithm in Figure 8.3c
   $\mathbf{W}_{\cdot h} = \mathbf{W} \text{diag}(\mathbf{b})$ 
  for  $g = 1 : p_{\dot{A}}$ 
    compute  $\mathbf{a} = \mathbf{A}(\cdot, g)$  using algorithm in Figure 8.3c
     $\mathbf{W}_{gh} = \text{diag}(\mathbf{a}) \mathbf{W}_{\cdot h}$ 
    compute  $\mathbf{Z} = \mathbf{A}'_s \mathbf{W}_{gh} \mathbf{B}_t$  using algorithms in Figure 9.1a or Figure 9.1b
     $(\mathbf{X}'_i \mathbf{W} \mathbf{X}_j)((gp_{\dot{A}_s} - p_{\dot{A}_s} + 1) : gp_{\dot{A}_s}, (hp_{\dot{B}_t} - p_{\dot{B}_t} + 1) : hp_{\dot{B}_t}) = \mathbf{Z}$ 

```

---

**Figure 9.4:** A concise presentation of the algorithms in Figure 9.2a and Figure 9.2b.  $\mathbf{X}_i$  is a tensor with  $s$  packed margins from  $(\bar{\mathbf{A}}_1, \mathbf{k}_{\bar{\mathbf{A}}_1}, m_{\bar{\mathbf{A}}_1}, p_{\bar{\mathbf{A}}_1})$  to  $(\bar{\mathbf{A}}_s, \mathbf{k}_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{A}}_s}, p_{\bar{\mathbf{A}}_s})$ , and  $\mathbf{X}_j$  is a tensor with  $t$  packed margins from  $(\bar{\mathbf{B}}_1, \mathbf{k}_{\bar{\mathbf{B}}_1}, m_{\bar{\mathbf{B}}_1}, p_{\bar{\mathbf{B}}_1})$  to  $(\bar{\mathbf{B}}_t, \mathbf{k}_{\bar{\mathbf{B}}_t}, m_{\bar{\mathbf{B}}_t}, p_{\bar{\mathbf{B}}_t})$ . It is legitimate to have  $s = 1$  and / or  $t = 1$  so that a tensor degenerates to a singleton.

### 9.1.3 A unified representation of $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$ computation

In fact, there is a unified representation for the computation of  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$ . Note that a singleton can be seen as a one-margin tensor, so I don't have to distinguish a singleton and a tensor when making discussion.

Let  $\mathbf{X}_i = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)$  and  $\mathbf{X}_j = (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)$ , and define  $\dot{\mathbf{A}} = \mathbf{1}$  (a single-column matrix of 1) and  $p_{\dot{A}} = 1$  for  $s = 1$ . Similarly define  $\dot{\mathbf{B}} = \mathbf{1}$  and  $p_{\dot{B}} = 1$  for  $t = 1$ . Then by applying (1.1) to both  $\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s$  and  $\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t$ , there is

$$(\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t) = \begin{bmatrix} \cdots & \mathbf{A}'_s \text{diag}(\dot{\mathbf{A}}(\cdot, g)) \mathbf{W} \text{diag}(\dot{\mathbf{B}}(\cdot, h)) \mathbf{B}_t & \cdots \end{bmatrix}.$$

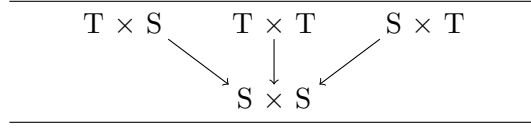
The RHS block matrix has  $p_{\dot{A}} \times p_{\dot{B}}$  block elements.

- If  $s > 1$  and  $t > 1$  (i.e., both  $\mathbf{X}_i$  and  $\mathbf{X}_j$  are tensors), the RHS block matrix has  $s$  block rows and  $t$  block columns. The weight matrix for the  $(g, h)^{\text{th}}$  block is  $\text{diag}(\dot{\mathbf{A}}(\cdot, g)) \mathbf{W} \text{diag}(\dot{\mathbf{B}}(\cdot, h))$ .
- If  $s > 1$  and  $t = 1$  (i.e.,  $\mathbf{X}_i$  is a tensor and  $\mathbf{X}_j$  is a singleton), the RHS block matrix has a block column of  $p_{\dot{A}}$  block elements. The weight matrix for the  $g^{\text{th}}$  block reduces to  $\text{diag}(\dot{\mathbf{A}}(\cdot, g)) \mathbf{W}$ .
- If  $s = 1$  and  $t > 1$  (i.e.,  $\mathbf{X}_i$  is a singleton and  $\mathbf{X}_j$  is a tensor), the RHS block matrix has a block row of  $p_{\dot{B}}$  block elements. The weight matrix for the  $h^{\text{th}}$  block reduces to  $\mathbf{W} \text{diag}(\dot{\mathbf{B}}(\cdot, h))$ .
- If  $s = t = 1$  (i.e., both  $\mathbf{X}_i$  and  $\mathbf{X}_j$  are singletons),  $p_{\dot{A}} = p_{\dot{B}} = 1$  hence the RHS block matrix only has one block element  $\mathbf{A}'_s \mathbf{W} \mathbf{B}_t$ ;

In this way, algorithms in Figure 9.2a and Figure 9.2b can be concisely presented by algorithm in Figure 9.4 (and there is also no need to do the matrix transpose in Figure 9.2b). The computational idea is now described by Figure 9.5. This idea is easier to programme in practice than that in Figure 9.5, since writing a single routine with a double loop nest is sufficient to handle all  $S \times T$ ,  $T \times S$  and  $T \times T$  cases. By contrast, the design in Figure 9.3 requires writing separate routines for these cases. However, in whichever design,  $\mathbf{X}'_i \mathbf{W} \mathbf{X}_j$  of  $S \times S$  type is the core.

## 9.2 FLOP comparison with previous discrete algorithms

So far the FLOP count for bamboos algorithms has only been given for  $S \times S$  (see Figure 9.1a and Figure 9.1b). I will now summarize FLOP count for all  $S \times S$ ,  $S \times T$ ,  $T \times S$  and  $T \times T$  blocks in



**Figure 9.5:** Illustrating the computational idea of  $\mathbf{X}_i' \mathbf{W} \mathbf{X}_j$  with a dependency graph, following algorithm in Figure 9.4. This design is easier to programme in practice than that in Figure 9.3 because a single routine with a double loop nest is sufficient to handle all  $S \times T$ ,  $T \times S$  and  $T \times T$  cases. By contrast, the sequential dependency in Figure 9.3 requires writing separate routines for these cases.

§9.2.1. Since **bamboos** is motivated to improve the previous discrete algorithms in Chapter 8, it is useful to first examine if there is any savings in FLOP count. §9.2.2 will restate the FLOP count for the old discrete algorithms, and comparison will be made in §9.2.3.

### 9.2.1 FLOP count for bamboos algorithms

Assuming  $\mathbf{A}_s$  and  $\mathbf{B}_t$  are two singletons, the following formally summarizes FLOP count for matrix cross-product of  $S \times S$  type.

- For diagonal  $\mathbf{W}$  (see algorithms in Figure 9.1a):
  - $\mathbf{A}_s' \mathbf{W} \mathbf{A}_s$ :  $\alpha_1 = O(n) + O(m_{\bar{A}_s} p_{\bar{A}_s}^2)$ ;
  - $\mathbf{A}_s' \mathbf{W} \mathbf{B}_t$ :  $\alpha_2 = O(n) + O(K_{st}^{[1]} p_{\bar{B}_t}) + O(m_{\bar{A}_s} p_{\bar{A}_s} p_{\bar{B}_t})$  for right-to-left computation and  $\alpha_3 = O(n) + O(K_{st}^{[1]} p_{\bar{A}_s}) + O(m_{\bar{B}_t} p_{\bar{A}_s} p_{\bar{B}_t})$  for left-to-right computation, where  $K_{st}^{[1]}$  is the number of non-zero elements in  $\bar{\mathbf{W}}$ , which is upper bounded by  $\min\{m_{\bar{A}_s} m_{\bar{B}_t}, n\}$ .
- For tri-diagonal  $\mathbf{W}$  (see algorithm Figure 9.1b; note the  $3n$  (or precisely  $(3n - 2)$ ) as there are three diagonals):
  - $\mathbf{A}_s' \mathbf{W} \mathbf{A}_s$ :  $\alpha_4 = O(n) + O(K_{ss}^{[3]} p_{\bar{A}_s}) + O(m_{\bar{A}_s} p_{\bar{A}_s}^2)$ , where  $K_{ss}^{[3]}$  is the number of non-zero elements in  $\bar{\mathbf{W}}$ , which is upper bounded by  $\min\{m_{\bar{A}_s}^2, 3n\}$ ;
  - $\mathbf{A}_s' \mathbf{W} \mathbf{B}_t$ :  $\alpha_5 = O(n) + O(K_{st}^{[3]} p_{\bar{B}_t}) + O(m_{\bar{A}_s} p_{\bar{A}_s} p_{\bar{B}_t})$  for right-to-left computation and  $\alpha_6 = O(n) + O(K_{st}^{[3]} p_{\bar{A}_s}) + O(m_{\bar{B}_t} p_{\bar{A}_s} p_{\bar{B}_t})$  for left-to-right computation, where  $K_{st}^{[3]}$  is the number of non-zero elements in  $\bar{\mathbf{W}}$ , which is upper bounded by  $\min\{m_{\bar{A}_s} m_{\bar{B}_t}, 3n\}$ .

Note that

- $K$  is indexed by  $s$  and  $t$  to imply its dependence on  $\mathbf{k}_{\bar{A}_s}$  and  $\mathbf{k}_{\bar{B}_t}$ . The superscript implies its dependence on the structure of  $\mathbf{W}$ : “[1]” for diagonal  $\mathbf{W}$  and “[3]” for tri-diagonal  $\mathbf{W}$ .
- A run-time choice between left-to-right computation and right-to-left computation can be made, so for example, the realistic FLOP count for  $\mathbf{A}_s' \mathbf{W} \mathbf{B}_t$  with diagonal  $\mathbf{W}$  is  $\min\{\alpha_2, \alpha_3\}$ .

The unified representation in Figure 9.1.3 (or algorithm in Figure 9.4) makes it straightforward to derive FLOP count for  $S \times T$ ,  $T \times S$  and  $T \times T$  blocks; see Table 9.1. Note that

- The computational complexity for weights aggregation in  $\alpha_1, \alpha_2, \dots, \alpha_6$  needs be adjusted to  $O(stn)$  to account for the computation of  $\mathbf{a}$  and  $\mathbf{b}$  in Figure 9.4.
- I have replaced the big “O” notation with exact FLOP count. The constant proportionality for weights aggregation is 1 for diagonal  $\mathbf{W}$  and 3 for tri-diagonal  $\mathbf{W}$ , while the constant proportionality for any matrix-matrix multiplication is 2.

**Table 9.1:** FLOP count for bamboos algorithms. The unified representation for  $\mathbf{X}_i$  and  $\mathbf{X}_j$  is used (see §9.1.3):  $\mathbf{X}_i$  is a tensor with  $s$  packed margins from  $(\bar{\mathbf{A}}_1, \mathbf{k}_{\bar{\mathbf{A}}_1}, m_{\bar{\mathbf{A}}_1}, p_{\bar{\mathbf{A}}_1})$  to  $(\bar{\mathbf{A}}_s, \mathbf{k}_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{A}}_s}, p_{\bar{\mathbf{A}}_s})$ , and  $\mathbf{X}_j$  is a tensor with  $t$  packed margins from  $(\bar{\mathbf{B}}_1, \mathbf{k}_{\bar{\mathbf{B}}_1}, m_{\bar{\mathbf{B}}_1}, p_{\bar{\mathbf{B}}_1})$  to  $(\bar{\mathbf{B}}_t, \mathbf{k}_{\bar{\mathbf{B}}_t}, m_{\bar{\mathbf{B}}_t}, p_{\bar{\mathbf{B}}_t})$ . It is legitimate to have  $s = 1$  and / or  $t = 1$  so that a tensor degenerates to a singleton, in which case there is  $p_{\bar{\mathbf{A}}} = 1$  and / or  $p_{\bar{\mathbf{B}}} = 1$ . Exact FLOP count instead of big “O” notation is given; see main text for more information.

$$\begin{aligned}
\alpha_1 &= stn + 2m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}_s}^2 \\
\alpha_2 &= stn + 2K_{st}^{[1]}p_{\bar{\mathbf{B}}_t} + 2m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{B}}_t}, \quad K_{st}^{[1]} \leq \min\{m_{\bar{\mathbf{A}}_s}m_{\bar{\mathbf{B}}_t}, n\} \\
\alpha_3 &= stn + 2K_{st}^{[1]}p_{\bar{\mathbf{A}}_s} + 2m_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{B}}_t}, \quad K_{st}^{[1]} \leq \min\{m_{\bar{\mathbf{A}}_s}m_{\bar{\mathbf{B}}_t}, n\} \\
\alpha_4 &= 3stn + 2K_{ss}^{[3]}p_{\bar{\mathbf{A}}_s} + 2m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}_s}^2, \quad K_{ss}^{[3]} \leq \min\{m_{\bar{\mathbf{A}}_s}^2, 3n\} \\
\alpha_5 &= 3stn + 2K_{st}^{[3]}p_{\bar{\mathbf{B}}_t} + 2m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{B}}_t}, \quad K_{st}^{[3]} \leq \min\{m_{\bar{\mathbf{A}}_s}m_{\bar{\mathbf{B}}_t}, 3n\} \\
\alpha_6 &= 3stn + 2K_{st}^{[3]}p_{\bar{\mathbf{A}}_s} + 2m_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{B}}_t}, \quad K_{st}^{[3]} \leq \min\{m_{\bar{\mathbf{A}}_s}m_{\bar{\mathbf{B}}_t}, 3n\}
\end{aligned}$$

	$\mathbf{X}_i' \mathbf{W} \mathbf{X}_i = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)$	$\mathbf{X}_i' \mathbf{W} \mathbf{X}_j = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)$
diagonal $\mathbf{W}$	$p_{\bar{\mathbf{A}}}^2 \alpha_1$	$p_{\bar{\mathbf{A}}} p_{\bar{\mathbf{B}}} \min\{\alpha_2, \alpha_3\}$
tri-diagonal $\mathbf{W}$	$p_{\bar{\mathbf{A}}}^2 \alpha_4$	$p_{\bar{\mathbf{A}}} p_{\bar{\mathbf{B}}} \min\{\alpha_5, \alpha_6\}$

## 9.2.2 A recap of FLOP count for the previous discrete algorithms

FLOP count for the old discrete algorithms was provided in captions of Figure 8.5a and Figure 8.5b, and the following is a recap.

- $\mathbf{A}_s' \mathbf{W} \mathbf{B}_t$ :  $O(np_{\bar{\mathbf{B}}_t}) + O(m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{B}}_t})$ ;
- $(\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} \mathbf{B}_t$ :  $O(np_{\bar{\mathbf{B}}_t}) + O(snp_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}/b) + O(np_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}) + O(m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{B}}_t})$ ;
- $\mathbf{A}_s' \mathbf{W} (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)$ :  $O(ntp_{\bar{\mathbf{B}}_t}) + O(m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t})$ ;
- $(\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)$ :  $O(ntp_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t}) + O(snp_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t}/b) + O(np_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t}) + O(m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t})$ .

If the column chunk size  $b$  is set to  $20s$  as is discussed in §8.4, the  $O(snp_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t}/b)$  for example, will be negligible compared with  $O(np_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t})$ . The unified representation between singleton and tensor can then be applied to simplify the results for all 4 cases to  $O(ntp_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t}) + O(np_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t}) + O(m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t})$ . To ease comparison with FLOP counts for bamboos algorithms (see Table 9.1), it is helpful to replace the big “O” notation with exact FLOP count. This proceeds as follows.

- The  $O(ntp_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t})$  complexity is associated with unpacking and reweighting of  $\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t$ . For the former, the constant proportionality is 1; for the latter, the constant proportionality is 1 if  $\mathbf{W}$  is diagonal and 6 if  $\mathbf{W}$  is tri-diagonal.
- The  $O(np_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t})$  complexity is for matrix row-aggregation. The constant proportionality is 1.
- The  $O(m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t})$  complexity is for matrix cross-product, so constant proportionality is 2.

To summarize, the FLOP count for computing  $(\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)$  is

- $2ntp_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t} + np_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t} + 2m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t}$  for diagonal  $\mathbf{W}$ ;
- $7ntp_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t} + np_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t} + 2m_{\bar{\mathbf{A}}_s}p_{\bar{\mathbf{A}}}p_{\bar{\mathbf{B}}_t}p_{\bar{\mathbf{B}}_t}$  for tri-diagonal  $\mathbf{W}$ ;

Still, recall that  $(\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)' \mathbf{W} (\dot{\mathbf{A}} \tilde{\otimes} \bar{\mathbf{A}}_s)$  may be computed instead of  $(\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)$  if it is cheaper. The FLOP count for this transposed variant is



**Table 9.2:** FLOP count for previous discrete algorithms in Chapter 8. Raw results are rearranged to the following form to facilitate comparison with Table 9.1; see main text for details of derivation.

$$\begin{aligned}
\gamma_1 &= np_{\bar{A}_s}(1 + 2s/p_{\bar{A}}) + 2m_{\bar{A}_s}p_{\bar{A}_s}^2 \\
\gamma_2 &= np_{\bar{B}_t}(1 + 2t/p_{\bar{A}}) + 2m_{\bar{A}_s}p_{\bar{A}_s}p_{\bar{B}_t} \\
\gamma_3 &= np_{\bar{A}_s}(1 + 2s/p_{\bar{B}}) + 2m_{\bar{B}_t}p_{\bar{A}_s}p_{\bar{B}_t} \\
\gamma_4 &= np_{\bar{A}_s}(1 + 7s/p_{\bar{A}}) + 2m_{\bar{A}_s}p_{\bar{A}_s}^2 \\
\gamma_5 &= np_{\bar{B}_t}(1 + 7t/p_{\bar{A}}) + 2m_{\bar{A}_s}p_{\bar{A}_s}p_{\bar{B}_t} \\
\gamma_6 &= np_{\bar{A}_s}(1 + 7s/p_{\bar{B}}) + 2m_{\bar{B}_t}p_{\bar{A}_s}p_{\bar{B}_t}
\end{aligned}$$

	$\mathbf{X}_i' \mathbf{W} \mathbf{X}_i = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)$	$\mathbf{X}_i' \mathbf{W} \mathbf{X}_j = (\dot{\mathbf{A}} \tilde{\otimes} \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{B}} \tilde{\otimes} \mathbf{B}_t)$
diagonal $\mathbf{W}$	$p_{\bar{A}}^2 \gamma_1$	$p_{\bar{A}} p_{\bar{B}} \min\{\gamma_2, \gamma_3\}$
tri-diagonal $\mathbf{W}$	$p_{\bar{A}}^2 \gamma_4$	$p_{\bar{A}} p_{\bar{B}} \min\{\gamma_5, \gamma_6\}$

- $2nsp_{\bar{A}}p_{\bar{A}_s} + np_{\bar{A}}p_{\bar{B}}p_{\bar{A}_s} + 2m_{\bar{B}_t}p_{\bar{B}}p_{\bar{B}_t}p_{\bar{A}}p_{\bar{A}_s}$  for diagonal  $\mathbf{W}$ ;
- $7nsp_{\bar{A}}p_{\bar{A}_s} + np_{\bar{A}}p_{\bar{B}}p_{\bar{A}_s} + 2m_{\bar{B}_t}p_{\bar{B}}p_{\bar{B}_t}p_{\bar{A}}p_{\bar{A}_s}$  for tri-diagonal  $\mathbf{W}$ ;

To facilitate comparison with **bamboos**, it helps to rearrange above results in a similar format to Table 9.1. This only involves a little more paper work, and Table 9.2 presents such results. In this way, it is sufficient to compare  $\alpha_i$  with  $\gamma_i$  to draw conclusions.

### 9.2.3 FLOP count comparison and summary

It is not easy to compare FLOP count of **bamboos** algorithms and the old discrete algorithms. A message from Table 9.1 is that realistic complexity of **bamboos** algorithms is data dependent via those  $K$  values. It is very crude to replace  $K$  values by their upper bound, but this at least provides some idea / guidance.

In the best case where  $m_{\bar{A}_s}m_{\bar{B}_t} < n$  or  $m_{\bar{A}_s}m_{\bar{B}_t} < 3n$ , there are

$$\begin{aligned}
\alpha_1 &= stn + 2m_{\bar{A}_s}p_{\bar{A}_s}^2, \\
\alpha_2 &\leq stn + 2m_{\bar{A}_s}m_{\bar{B}_t}p_{\bar{B}_t} + 2m_{\bar{A}_s}p_{\bar{A}_s}p_{\bar{B}_t}, \\
\alpha_3 &\leq stn + 2m_{\bar{A}_s}m_{\bar{B}_t}p_{\bar{A}_s} + 2m_{\bar{B}_t}p_{\bar{A}_s}p_{\bar{B}_t}, \\
\alpha_4 &\leq 3stn + 2m_{\bar{A}_s}^2p_{\bar{A}_s} + 2m_{\bar{A}_s}p_{\bar{A}_s}^2, \\
\alpha_5 &\leq 3stn + 2m_{\bar{A}_s}m_{\bar{B}_t}p_{\bar{B}_t} + 2m_{\bar{A}_s}p_{\bar{A}_s}p_{\bar{B}_t}, \\
\alpha_6 &\leq 3stn + 2m_{\bar{A}_s}m_{\bar{B}_t}p_{\bar{A}_s} + 2m_{\bar{B}_t}p_{\bar{A}_s}p_{\bar{B}_t},
\end{aligned}$$

so that the only obviously  $n$  dependent terms are associated with weights aggregation. In large GAM fitting, it is common that 1)  $n \rightarrow \infty$ , 2)  $m_{\bar{A}_s}$  and  $m_{\bar{B}_t}$  grows sublinearly with  $n$  and 3)  $p_{\bar{A}_s}$  and  $p_{\bar{B}_t}$  grows sublinearly with  $m_{\bar{A}_s}$  and  $m_{\bar{B}_t}$ . Note that since there is  $m_{\bar{A}_s}m_{\bar{B}_t} < n$  or  $m_{\bar{A}_s}m_{\bar{B}_t} < 3n$ , the growth rate of  $m_{\bar{A}_s}$  and  $m_{\bar{B}_t}$  is at most  $n^{\frac{1}{2}}$ . Let us consider two scenarios.

- $n \rightarrow \infty$  and the growth rate of  $m_{\bar{A}_s}$  and  $m_{\bar{B}_t}$  is smaller than  $n^{\frac{1}{2}}$ . In this case, complexity in weights aggregation is dominant. Asymptotically there are  $\alpha_1 \approx \alpha_2 \approx \alpha_3 \approx stn$ ,  $\alpha_4 \approx \alpha_5 \approx \alpha_6 \approx 3stn$ ,  $\gamma_1 \approx \gamma_3 \approx \gamma_4 \approx \gamma_6 \approx np_{\bar{A}_s}$  and  $\gamma_2 \approx \gamma_5 \approx np_{\bar{B}_t}$ . **bamboos** algorithms yield huge computational savings.
- $n \rightarrow \infty$  and the growth rate of  $m_{\bar{A}_s}$  and  $m_{\bar{B}_t}$  attains  $n^{\frac{1}{2}}$ . In this case, complexity in matrix-matrix multiplication accounts for a fixed proportion of gross complexity. It may also be non-negligible compared with weights aggregation. Analysis based on upper bound gives no useful conclusion, because for example, if  $m_{\bar{A}_s}m_{\bar{B}_t}$  is almost  $n$  for diagonal  $\mathbf{W}$ , the upper bound of  $\alpha_3$  is almost

$2np_{\bar{A}_s}$ , readily larger than  $\gamma_3 \approx np_{\bar{A}_s}$ . If  $m_{\bar{A}_s}m_{\bar{B}_t}$  is almost  $3n$  for tri-diagonal  $\mathbf{W}$ , the upper bound of  $\alpha_6$  is almost  $6np_{\bar{A}_s}$ , even larger than  $\gamma_6 \approx np_{\bar{A}_s}$ . However, if  $m_{\bar{A}_s}m_{\bar{B}_t}$  is really this large, the real situation is likely to be that  $\bar{\mathbf{W}}$  is very sparse, therefore the upper bound used here is too much bigger than  $K$  values that actually determine the computational complexity.

For the worst case where  $m_{\bar{A}_s}m_{\bar{B}_t} \geq n$  or  $m_{\bar{A}_s}m_{\bar{B}_t} \geq 3n$ , the situation is basically as same as the second limiting case above. No precise comparison can be made when the sparsity of  $\bar{\mathbf{W}}$  steps in. However, this does not imply that **bamboos** has no advantage. In fact, the advantage may be magnificent. A sparse  $\bar{\mathbf{W}}$  can mean that  $K \approx 0.1n$ ,  $K \approx 0.01n$  or even  $K \approx 0.001n$ , therefore the first limiting case above is asymptotically recovered.

To some degree, the superiority of **bamboos**, if any, can only be established on practical benchmarking. This will be done later in §9.4. However, at least one firm conclusion can be made from the analysis so far: it is important to make weights aggregation as fast as possible, as that is what **bamboos** algorithms will be primarily doing in the  $n \rightarrow \infty$  limiting case. In the next section, I will demonstrate how to maximize the performance of weights aggregation. In addition, if you have read previous sections carefully, you may have already realized that although I have mentioned sparsity many times by now, nothing is really said about how to detect whether  $\bar{\mathbf{W}}$  is dense or sparse, and if sparse, how to represent it. The reason is that the “dense-sparse switch” for  $\bar{\mathbf{W}}$  construction as well as the sparse representation of  $\bar{\mathbf{W}}$  is embedded in the performance maximization of weights aggregation in **bamboos** implementation. This will be clear in the next section.

### 9.3 Fast weights aggregation for $\bar{\mathbf{W}}$ and its sparse construction

The way to maximize the performance of weights aggregation is subtle. Since it has been established that computing  $\mathbf{X}_i' \mathbf{W} \mathbf{X}_j$  of  $S \times S$  type plays a pivotal role in **bamboos** (see Figure 9.5), it might appear that performance optimization of weights aggregation in computing  $S \times S$  blocks is the way to go and any gains will automatically be broadcasted to computations of other types of blocks. However, this is not true. In fact, there is no way to enhance the performance of weights aggregation in  $S \times S$  case; optimization opportunity is found instead in the repeated application of the same kind of weights aggregation in  $S \times T$ ,  $T \times S$  and  $T \times T$  blocks. In the rest of this section,

- §9.3.1 introduces the idea of nested sorting, run-length encoding and bin aggregation that makes weights aggregation as fast as possible;
- §9.3.2 explains the sorting algorithms used in the nested sorting;
- §9.3.3 demonstrates how nested sorting, run-length encoding and bin aggregation facilitates sparse representation of  $\bar{\mathbf{W}}$ , and how the “dense-sparse switch” for  $\bar{\mathbf{W}}$  works in **bamboos**.

Note that for clarity, all these sections are discussed using diagonal  $\mathbf{W}$ . How tri-diagonal  $\mathbf{W}$  is dealt with will be briefly mentioned in §9.3.4 in the end.

#### 9.3.1 Nested sorting, run-length encoding and bin aggregation

Consider constructing  $\bar{\mathbf{W}}$  as a dense matrix when computing an  $S \times S$  block  $\mathbf{A}_s' \mathbf{W} \mathbf{B}_t$  where  $\mathbf{W}$  is diagonal. Algorithm in the right panel of Figure 9.1a hints that elements of  $\bar{\mathbf{W}}$  are likely to be accessed randomly over the entire matrix during the , as entries of index vectors  $\mathbf{k}_{\bar{A}}$  and especially  $\mathbf{k}_{\bar{B}}$  are not sorted. Therefore, the  $O(n)$  random memory access behind the  $O(n)$  FLOP in aggregation

original pair	primary sort		secondary sort	
$(\mathbf{k}_{\bar{A}_s}, \mathbf{k}_{\bar{B}_t})$	$(\mathbf{k}_{\bar{A}_s}^{[1]}, \mathbf{k}_{\bar{B}_t}^{[1]})$	$\sigma^{[1]}$	$(\mathbf{k}_{\bar{A}_s}^{[2]}, \mathbf{k}_{\bar{B}_t}^{[2]})$	$\sigma^{[2]}$
(1, 4)	(2, 1)	5	(2, 1)	5
(4, 4)	(3, 1)	9	(2, 1)	12
(1, 3)	(2, 1)	12	(3, 1)	9
(4, 3)	(1, 2)	10	(1, 2)	10
(2, 1)	(2, 2)	11	(2, 2)	11
(4, 3)	(1, 3)	3	(1, 3)	3
(1, 4)	(4, 3)	4	(1, 3)	16
(3, 3)	(4, 3)	6	(3, 3)	8
(3, 1)	(3, 3)	8	(4, 3)	4
(1, 2)	(1, 3)	16	(4, 3)	6
(2, 2)	(1, 4)	1	(1, 4)	1
(2, 1)	(4, 4)	2	(1, 4)	7
(4, 4)	(1, 4)	7	(1, 4)	15
(4, 4)	(4, 4)	13	(4, 4)	2
(1, 4)	(4, 4)	14	(4, 4)	13
(1, 3)	(1, 4)	15	(4, 4)	14

**Figure 9.6:** An illustration of the nested sorting on  $(\mathbf{k}_{\bar{A}_s}, \mathbf{k}_{\bar{B}_t})$  pairs, using an example with  $m_{\bar{A}_s} = m_{\bar{B}_t} = 4$  and  $n = 16$ . Primary sort arranges  $\mathbf{k}_{\bar{B}_t}$  into  $\mathbf{k}_{\bar{B}_t}^{[1]}$  whose entries are in non-descending order, while generating a permutation index vector  $\sigma^{[1]}$ .  $\mathbf{k}_{\bar{A}_s}$  is permuted accordingly to  $\mathbf{k}_{\bar{A}_s}^{[1]}$ . At this stage  $\mathbf{k}_{\bar{B}_t}^{[1]}$  can be cut into several “bins”, i.e., segments with runs of the value. Secondary sort is a bin-wise sort. Entries of  $\mathbf{k}_{\bar{A}_s}^{[1]}$  in each bin are sorted in non-descending order to  $\mathbf{k}_{\bar{A}_s}^{[2]}$  and  $\sigma^{[1]}$  is permuted accordingly to  $\sigma^{[2]}$ . Note that this stage has no effect on  $\mathbf{k}_{\bar{B}_t}^{[1]}$  so  $\mathbf{k}_{\bar{B}_t}^{[2]} = \mathbf{k}_{\bar{B}_t}^{[1]}$ . In the end of the nested sorting,  $(\mathbf{k}_{\bar{A}_s}^{[2]}, \mathbf{k}_{\bar{B}_t}^{[2]})$  admits a bin structure.

is generally much more expensive, when  $\bar{\mathbf{W}}$  is not small enough to fit in CPU cache. For example, a  $100 \times 100$  matrix of double-precision floating-point numbers is already 150% larger than a 32 KB L1 data cache.

Such inefficiency can be tolerated in computing the above  $S \times S$  block. But consider, for example the computation of a  $T \times T$  block  $(\dot{\mathbf{A}} \otimes \mathbf{A}_s)' \mathbf{W} (\dot{\mathbf{B}} \otimes \mathbf{B}_t)$  for diagonal  $\mathbf{W}$ , the same type of weights aggregation is repeatedly performed  $p_{\dot{\mathbf{A}}} p_{\dot{\mathbf{B}}}$  times through the iteration computation of  $\mathbf{A}_s' \mathbf{W}_{gh} \mathbf{B}_t$  (see Figure 9.4). By “the same type” I mean that the access pattern to  $\bar{\mathbf{W}}_{gh}$  is exactly the same for any  $g$  and  $h$ ; only values of  $\mathbf{W}_{gh}$  and  $\bar{\mathbf{W}}_{gh}$  are changing from iteration to iteration. If this common access pattern can be “memorized” somehow, i.e., the loopup to  $\mathbf{k}_{\bar{A}_s}$  and  $\mathbf{k}_{\bar{B}_t}$  is only done once rather than  $p_{\dot{\mathbf{A}}} p_{\dot{\mathbf{B}}}$  times, the efficiency of weights aggregation would be much improved.

One way to do this is by a nested sorting on  $(\mathbf{k}_{\bar{A}_s}, \mathbf{k}_{\bar{B}_t})$ , which includes the following two stages.

- A *primary sort*  $\sigma^{[1]} : (\mathbf{k}_{\bar{A}_s}, \mathbf{k}_{\bar{B}_t}) \rightarrow (\mathbf{k}_{\bar{A}_s}^{[1]}, \mathbf{k}_{\bar{B}_t}^{[1]})$ , that sorts  $\mathbf{k}_{\bar{B}_t}$  in non-descending order to  $\mathbf{k}_{\bar{B}_t}^{[1]}$ . Note that  $\mathbf{k}_{\bar{A}_s}$  needs be permuted according to the permutation index vector  $\sigma^{[1]}$  to  $\mathbf{k}_{\bar{A}_s}^{[1]}$ , so that  $(\mathbf{k}_{\bar{A}_s}^{[1]}, \mathbf{k}_{\bar{B}_t}^{[1]})$  and  $(\mathbf{k}_{\bar{A}_s}, \mathbf{k}_{\bar{B}_t})$  have the same pairwise relationship.
- A *secondary sort*  $\sigma^* : (\mathbf{k}_{\bar{A}_s}^{[1]}, \mathbf{k}_{\bar{B}_t}^{[1]}, \sigma^{[1]}) \rightarrow (\mathbf{k}_{\bar{A}_s}^{[2]}, \mathbf{k}_{\bar{B}_t}^{[2]}, \sigma^{[2]})$ , that further sorts entries of  $\mathbf{k}_{\bar{A}_s}^{[1]}$  in non-descending order to  $\mathbf{k}_{\bar{A}_s}^{[2]}$  for each “bin” (i.e., segment of the same value) of the sorted  $\mathbf{k}_{\bar{B}_t}^{[1]}$ . The permutation index  $\sigma^*$  is not of interest; rather,  $\sigma^{[2]}$  should be recorded. Note that this bin-wise sorting has no effect on  $\mathbf{k}_{\bar{B}_t}^{[1]}$  so  $\mathbf{k}_{\bar{B}_t}^{[2]} = \mathbf{k}_{\bar{B}_t}^{[1]}$ .

At the end of the nested sorting,  $(\mathbf{k}_{\bar{A}_s}^{[2]}, \mathbf{k}_{\bar{B}_t}^{[2]})$  will be partitioned into “bins”. This ensures that during weights aggregation, elements of  $\bar{\mathbf{W}}$  will be accessed sequentially. See Figure 9.6 for an illustration.

The nested sorting changes the way that  $\bar{\mathbf{W}}_{gh}$  is computed. The resulting permutation  $\sigma^{[2]}$  must be applied to  $\mathbf{W}_{gh} = \text{diag}(\mathbf{w}_{gh})$ . An implication is that while elements of  $\bar{\mathbf{W}}_{gh}$  are accessed sequentially, elements of  $\mathbf{w}_{gh}$  are now accessed in a random fashion. This reverses the situation before nested

sorting where elements of  $\bar{\mathbf{W}}_{gh}$  are accessed randomly and elements of  $\mathbf{w}_{gh}$  are accessed sequentially. No practical gains will be attained in this way. Fortunately, there is a workaround. Note that the permutation of  $\mathbf{w}_{gh}$  is

$$\dot{\mathbf{w}}_{gh} := \mathbf{w}_{gh}(\boldsymbol{\sigma}^{[2]}) = \text{diag}(\dot{\mathbf{A}}(\boldsymbol{\sigma}^{[2]}, g))\text{diag}(\mathbf{w}(\boldsymbol{\sigma}^{[2]}))\text{diag}(\dot{\mathbf{B}}(\boldsymbol{\sigma}^{[2]}, h)) = \dot{\mathbf{A}}(\boldsymbol{\sigma}^{[2]}, g) \circ \mathbf{w}(\boldsymbol{\sigma}^{[2]}) \circ \dot{\mathbf{B}}(\boldsymbol{\sigma}^{[2]}, h),$$

where “ $\circ$ ” is the Hadamard product between two vectors (see the footnote in §8.3 if you forget what it is). The following observations and solutions can then be made.

- $\mathbf{w}(\boldsymbol{\sigma}^{[2]})$  is loop-invariant w.r.t.  $g$  and  $h$  so it can be pre-computed.
- Computation of  $\dot{\mathbf{A}}(\boldsymbol{\sigma}^{[2]}, g)$  and  $\dot{\mathbf{B}}(\boldsymbol{\sigma}^{[2]}, h)$  is basically no different to computation of  $\dot{\mathbf{A}}(\cdot, g)$  and  $\dot{\mathbf{B}}(\cdot, h)$ . Algorithm in Figure 8.3c still applies, by now using  $\mathbf{k}_{\bar{A}_u}(\boldsymbol{\sigma}^{[2]})$ ,  $u = 1, 2, \dots, (s-1)$  and  $\mathbf{k}_{\bar{B}_v}(\boldsymbol{\sigma}^{[2]})$ ,  $v = 1, 2, \dots, (v-1)$  as index vectors when unpacking columns of  $\bar{\mathbf{A}}_u$  and  $\bar{\mathbf{B}}_v$ . All these index vectors are loop-invariant hence can be pre-computed.

In this way, the negative effect of permutation is eliminated in the double loop nest. In practice, pre-permutation of  $\mathbf{w}$ ,  $\mathbf{k}_{\bar{A}_u}$ ,  $u = 1, 2, \dots, (s-1)$  and  $\mathbf{k}_{\bar{B}_v}$ ,  $v = 1, 2, \dots, (v-1)$  can be implemented in secondary sort as soon as  $\boldsymbol{\sigma}^{[2]}$  is obtained.

What the nestes sorting achieves in practice, is that it casts the aggregation of  $\mathbf{w}_{gh}$  over  $(\mathbf{k}_{\bar{A}_s}, \mathbf{k}_{\bar{B}_t})$  to an equivalent aggregation of  $\dot{\mathbf{w}}_{gh}$  over  $(\mathbf{k}_{\bar{A}_s}^{[2]}, \mathbf{k}_{\bar{B}_t}^{[2]})$ . While cache-unfriendly random memory access is a performance penalty for the former, sequential memory access is almost guaranteed for latter. I used “almost” because unpacking columns of  $\bar{\mathbf{A}}_u$  and  $\bar{\mathbf{B}}_v$  still involves random memory access. But this is what has to be paid for using packed storage and it exists in both aggregation problems. Note that since  $(\mathbf{k}_{\bar{A}_s}^{[2]}, \mathbf{k}_{\bar{B}_t}^{[2]})$  admits “bin” structures, aggregation of  $\dot{\mathbf{w}}_{gh}$  is a bin-wise aggregation. To facilitate such aggregation, perform a *run-length encoding* for  $(\mathbf{k}_{\bar{A}_s}^{[2]}, \mathbf{k}_{\bar{B}_t}^{[2]})$ , producing the following triplets:

- $\mathbf{i} = (2, 3, 1, 2, 1, 3, 4, 1, 4);$
- $\mathbf{j} = (1, 1, 2, 2, 3, 3, 3, 4, 4);$
- $\mathbf{f} = (2, 1, 1, 1, 2, 1, 2, 3, 3),$

which imply that  $\bar{\mathbf{W}}(\mathbf{i}(v), \mathbf{j}(v))$ ,  $v = 1, 2, \dots, 9$  is the sum of  $\mathbf{f}(v)$  consecutive elements in  $\dot{\mathbf{w}}_{gh}$ . The number of elements in those triplets is the number of bins (here, 9). It is also  $K_{st}^{[1]}$ , the number of non-zero elements in  $\bar{\mathbf{W}}$ . Given these triplets, the bin aggregation can be implemented with the double loop nest in Figure 9.7.

Note that while this seems more complicated than the single aggregation loop based on  $(\mathbf{k}_{\bar{A}_s}, \mathbf{k}_{\bar{B}_t})$  and  $\mathbf{w}_{gh}$  (see the right panel of Figure 9.1a), it is in practice faster, because of the optimized memory access and the use of accumulator variable  $a$  to reduce the amount of write-back. Figure 9.8 benchmarks two aggregation methods on Intel Core i5-2557M. Note that for fixed  $n$ , as  $\bar{\mathbf{W}}$  gets larger and larger, the simple aggregation method is increasingly time consuming, which reflects the performance penalty from random memory access. On the other hand, the operation time for bin aggregation almost stays constant. On a  $141 \times 141$   $\bar{\mathbf{W}}$  (with about 20000 elements), the bin aggregation method is 3.9 times as fast.

### 9.3.2 Sorting algorithms in bamboos

In this section, I will explain the sorting algorithms used for primary sort and secondary sort. There are tons of sorting methods in general, so strictly speaking, careful choice should be made between

---

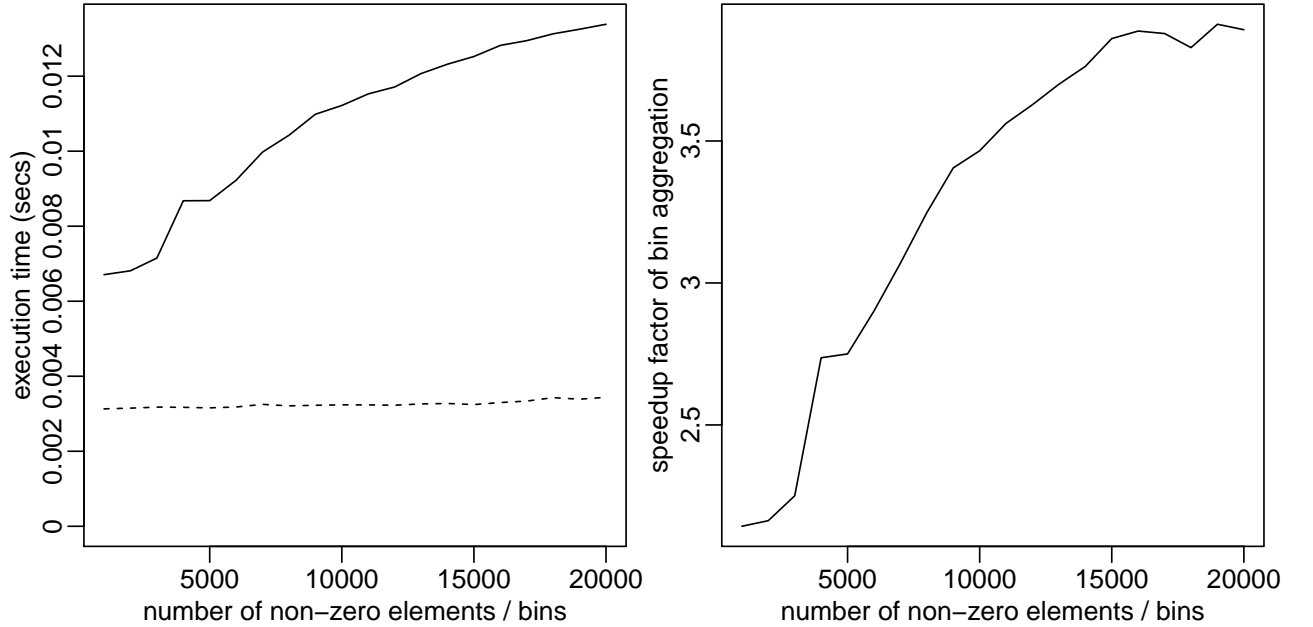
```

initialize  $\bar{\mathbf{W}}$  as a dense matrix
 $\bar{\mathbf{W}} = \text{zeros}(m_{\bar{\mathbf{A}}_s}, m_{\bar{\mathbf{B}}_t})$ 
number of elements to be skipped to reach the first element in a bin
 $s = 0$ 
loop through bins
for  $v = 1 : K_{st}^{[1]}$ 
    number of values to be aggregated
     $f = \mathbf{f}(v)$ 
    an accumulator variable in a CPU register
     $a = 0$ 
    aggregation on the  $v^{\text{th}}$  bin
    for  $u = 1 : f$ 
         $a += \dot{\mathbf{w}}_{gh}(s + u)$ 
    write-back
     $\bar{\mathbf{W}}(\mathbf{i}(v), \mathbf{j}(v)) = a$ 
     $s += f$ 

```

---

**Figure 9.7:** Fast bin aggregation that produces a dense  $\bar{\mathbf{W}}$ . The triplets  $\mathbf{i}$ ,  $\mathbf{j}$  and  $\mathbf{f}$  are the result of run-length encoding on  $(\mathbf{k}_{\bar{\mathbf{A}}_s}^{[2]}, \mathbf{k}_{\bar{\mathbf{B}}_t}^{[2]})$  (see main text for more information). Run-length encoding also reports the number of bins, which is in fact  $K_{st}^{[1]}$ , the number of non-zero elements in  $\bar{\mathbf{W}}$ . Note that talking about non-zero elements of  $\bar{\mathbf{W}}$  does not necessarily mean treating the matrix as sparse. However, the nested sorting, run-length encoding and bin aggregation do make it straightforward to construct  $\bar{\mathbf{W}}$  as sparse instead of dense. This will be clear in §9.3.3.



**Figure 9.8:** Benchmarking simple aggregation (see the *for* loop in the right panel of Figure 9.1a) and bin aggregation (see Figure 9.7) on Intel Core i5-2557M (see Appendix A for hardware information). A weight vector of length  $n = 10^6$  is fixed for the experiment, and the size of  $\bar{\mathbf{W}}$  grows. The left panel sketches the execution time (in seconds) for both methods (solid line for simple aggregation and dashed line for bin aggregation), and the right panel shows how faster bin aggregation is compared with simple aggregation. As  $\bar{\mathbf{W}}$  gets larger and larger, the simple aggregation method is increasingly time consuming, reflecting the performance penalty from random memory access. On the other hand, the operation time for bin aggregation almost stays constant. On a  $141 \times 141$   $\bar{\mathbf{W}}$  (with about 20000 elements), the bin aggregation method is 3.9 times as fast.

different methods. However, sorting is not a computational hotspot in **bamboos**, so there is no need to do premature optimization on this topic. Besides, sorting tasks in **bamboos** are integer sorting problems, and there are a few algorithms known to be best at such tasks. These include counting sort, radix sort (and maybe quicksort). At the moment, **bamboos** relies on counting sort for primary sort; this is guaranteed to be optimal. Binary MSD (most significant digit) radix sort is used for secondary sort; it may not be optimal but it is also good enough.

### Counting sort algorithm for primary sort

Primary sort is an integer sorting problem, where a length- $n$  vector  $\mathbf{k}_{\bar{B}_i}$  taking values in  $1, 2, \dots, m_{\bar{B}_i}$  needs be sorted in ascending order. In particular,  $n \gg m_{\bar{B}_i}$ . Counting sort is optimal for this task, with  $O(n)$  time complexity and  $O(m_{\bar{B}_i})$  space complexity. In the end of the sorting, it returns the required permutation index  $\sigma^{[1]}$  and a frequency table  $\mathbf{f}^{[1]}$  for  $1, 2, \dots, m$ . Note that  $\mathbf{k}_{\bar{B}_i}^{[1]}$  is not returned, but  $\mathbf{f}^{[1]}$  is equally informative on how  $(\mathbf{k}_{A_s}^{[1]}, \mathbf{k}_{\bar{B}_i}^{[1]})$  is cut into bins according to  $\mathbf{k}_{\bar{B}_i}^{[1]}$ .

The core idea of counting sort is as what its name suggests: counting. Consider the following length-20 toy integer vectors taking values  $1, 2, \dots, 7$ .

$$\mathbf{k}_{\bar{B}_i} = 3, 7, 2, 5, 1, 4, 6, 2, 5, 6, 1, 7, 3, 3, 4, 2, 1, 5, 6, 6$$

The first round of counting produces  $\mathbf{f}^{[1]}$ , and a vector of pointers  $\mathbf{p}$ :

	1	2	3	4	5	6	7
$\mathbf{f}^{[1]}$	3	3	3	2	3	4	2
$\mathbf{p}$	0	3	6	9	11	14	18

Here, a “header” is written in gray to help you associate values of  $\mathbf{f}^{[1]}$  and  $\mathbf{p}$  to unique integer values in  $\mathbf{k}_{\bar{B}_i}$ . If the goal of counting sort is just to obtain  $\mathbf{k}_{\bar{B}_i}^{[1]}$ , then the sorting readily completes, because the following sorted array is simply generated by replicating  $1, 2, \dots, 7$  according to the frequency table  $\mathbf{f}^{[1]}$ :

$$\mathbf{k}_{\bar{B}_i}^{[1]} = |1, 1, 1, |2, 2, 2, |3, 3, 3, |4, 4, |5, 5, 5, |6, 6, 6, 6, |7, 7$$

This also makes it easy to explain what  $\mathbf{p}$  is. It is just the pointer to those vertical bars: the number of elements before these bars.

When the goal of counting sort is to obtain  $\sigma^{[1]}$ , a bit more work is needed. We have to scan  $\mathbf{k}_{\bar{B}_i}$  for the second time to determine for example, which “1” in  $\mathbf{k}_{\bar{B}_i}$  is the first “1” in  $\mathbf{k}_{\bar{B}_i}^{[1]}$ , which “1” in  $\mathbf{k}_{\bar{B}_i}$  is the second “1” in  $\mathbf{k}_{\bar{B}_i}^{[1]}$ . This is essentially a second round of counting. To start with, a counter vector  $\mathbf{c}$  is initialized:

	1	2	3	4	5	6	7
$\mathbf{c}$	0	0	0	0	0	0	0

Then the scanning of  $\mathbf{k}_{\bar{B}_i}$  starts. The first value is  $\mathbf{k}_{\bar{B}_i}(1) = 3$ , so the counter  $\mathbf{c}(3)$  is incremented to record that the first “3” in  $\mathbf{k}_{\bar{B}_i}^{[1]}$  has been met. At the same time, it is easy to see that this “3” should go to position  $\mathbf{p}(3) + \mathbf{c}(3) = 7$  of  $\mathbf{k}_{\bar{B}_i}^{[1]}$ . In other words,  $\mathbf{k}_{\bar{B}_i}(1)$  is mapped to  $\mathbf{k}_{\bar{B}_i}^{[1]}(7)$  by the permutation, hence  $\sigma^{[1]}(7) = 1$ . The following *for* loop demonstrates how this works in general. Given this  $\sigma^{[1]}$ , permutating  $\mathbf{k}_{A_s}$  for  $\mathbf{k}_{A_s}^{[1]}$  is trivial.

second scanning of  $\mathbf{k}_{\bar{B}_i}$

for  $i = 1 : n$

  read in the  $i^{\text{th}}$  element

$v = \mathbf{k}_{\bar{B}_i}(i)$

  increment counter

$\mathbf{c}(v) += 1$

  write  $i$  to element  $(\mathbf{p}(v) + \mathbf{c}(v))$  of  $\sigma^{[1]}$

$\sigma^{[1]}(\mathbf{p}(v) + \mathbf{c}(v)) = i$

## Binary MSD radix sort for secondary sort

In secondary sort,  $\mathbf{k}_{A_s}^{[1]}$  is a length- $n$  integer vector taking values  $1, 2, \dots, m_{A_s}$  and there is  $n \gg m_{A_s}$ . The vector is then cut into  $m_{B_i}$  bins and bin-wise sort is applied on each bin to sort each segment in ascending order. Suppose the  $i^{\text{th}}$  segment of  $\mathbf{k}_{A_s}^{[1]}$  has length  $n_i$  and takes positive integer values no larger than  $m_i$ , then it is possible that  $m_i > n_i$  or even  $m_i \gg n_i$ . As a result, counting sort with  $O(n_i)$  time complexity and  $O(m_i)$  space complexity is very inefficient. It is a better idea to resort to other sorting methods with  $O(n_i \log(n_i))$  complexity on average, say quicksort and radix sort. Binary MSD radix sort is a simple variant of these algorithms. In the following I will demonstrate how this method works via a toy example.

Consider a length-15 integer vector below and its bit representation:

6,	2,	6,	3,	3,	3,	5,	1,	2,	3,	6,	6,	5,	3,	2
110,	010,	110,	011,	011,	011,	101,	001,	010,	011,	110,	110,	101,	011,	010

The maximum value in the vector is 6, so  $\lceil \log_2(6+1) \rceil = 3$  bits are sufficient for bit representation and the most significant digit is 3. The sorting algorithm first examines the 3rd digit (counting from right to left) of those binaries, and sorts them by this digit. In the following, this digit is displayed in black, while the rest of the digits are displayed in gray. Note the way that sorting is done: the vector is scanned from two sides in opposing directions. From left to right, the first “1” encountered is displayed in red; from right to left, the first “0” encountered is displayed in blue. The corresponding two numbers are swapped if the red digit “1” is before (or to the left of) the blue digit “0”. This process goes on until the two flagged digits cross each other somewhere in the middle. In this example, three swaps are made to reach this point.

110	010	110	011	011	011	101	001	010	011	110	110	101	011	010
010	010	110	011	011	011	101	001	010	011	110	110	101	011	110
010	010	011	011	011	011	101	001	010	011	110	110	101	110	110
010	010	011	011	011	011	011	001	010	101	110	110	101	110	110

Now the vector can be broken into two pieces at this crossover point: all numbers in the left piece have their 3rd digit being 0 and all numbers in the right piece have their 3rd digit being 1. The sorting algorithm then examines the 2nd digit on both pieces separately. In this example, it turns out that after one swap, a crossover point is hit in both pieces.

110	010	011	011	011	011	011	001	010	101	110	110	101	110	110
001	110	011	011	011	011	011	010	010	101	110	110	110	110	110

Now each piece can be further split into two, giving 4 pieces in total. The algorithm now examines the 1st digit on each piece separately. The first piece only has one single number hence is readily sorted; nothing needs be done further. The third and the fourth pieces are also readily sorted. Only the second piece needs further sorting.

001	010	011	011	011	011	011	010	010	101	101	110	110	110	110
001	010	010	011	011	011	011	010	011	101	101	110	110	110	110
001	010	010	010	011	011	011	011	011	101	101	110	110	110	110

The 1st digit is already the rightmost (least significant) digit, so the algorithm will terminate. The whole vector is now sorted, and I will now convert bit representation back to integers.

001,	010,	010,	010,	011,	011,	011,	011,	011,	101,	101,	110,	110,	110,	110
1,	2,	2,	2,	3,	3,	3,	3,	3,	5,	5,	6,	6,	6,	6

Here are a few remarks on this binary MSD radix sort.

- The algorithm is named with “binary” not just because it examines data in their binary format, but also because it scans the vector from two sides and eventually subdivides the vector into two. In this way, the algorithm can naturally be implemented recursively.

- The algorithm is an “in-place” algorithm. In the end of the algorithm, the initially unsorted vector will be overwritten by the sorted one.
- Permutation index vector can be generated alongside the sorted vector. This only requires that (in this example of 15 data) an index vector 1, 2, ..., 15 be initialized, and its elements swapped in the same way as elements of the vector are swapped. In fact, this can be more general. For example, in secondary sort,  $\sigma^{[1]}$  can be used instead of 1, 2, ...,  $n$ , and  $\sigma^{[2]}$  will overwrites  $\sigma^{[1]}$  when the sorting completes.

### 9.3.3 Dealing with sparse $\bar{W}$

Now I am going to address another key issue in **bamboos**: dealing with the possible sparsity in  $\bar{W}$ . Three questions need be answered:

1. How to construct / represent / store  $\bar{W}$  as a sparse matrix?
2. How to do matrix-matrix multiplication that involves a sparse  $\bar{W}$ ?
3. When is using a sparse  $\bar{W}$  advantageous to using a dense  $\bar{W}$ ?

Answering the first two questions just requires some basic knowledge of sparse matrices. In particular, neither question is performance related, as for example, there is no problem in representing / storing a realistic dense matrix using sparse matrix representation / storage method. The matter is whether doing so is worthwhile, which is exactly the aim of the last question.

#### Compressed column storage

There are undoubtedly many storage format for a sparse matrix. For a general sparse matrix without any special sparse structure<sup>1</sup>, common storage format include triplet storage format, compressed row storage (CRS) format and compressed column storage (CCS) format. There is no “right” storage format, but only the convenient format that fits into the application context. For example, the run-length encoding and bin aggregation following the nested sorting readily yields the triplet storage format. Consider the example in Figure 9.6, this is

- $\mathbf{i} = (2, 3, 1, 2, 1, 3, 4, 1, 4);$
- $\mathbf{j} = (1, 1, 2, 2, 3, 3, 3, 4, 4);$
- $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9),$

implying that  $\bar{W}(\mathbf{i}(v), \mathbf{j}(v)) = \mathbf{x}(v)$ ,  $v = 1, 2, \dots, 9$  where  $\mathbf{x}(v)$  is the result of bin aggregation. It is also straightforward to obtain the CCS format by replacing  $\mathbf{j}$  with the cumulative number of non-zero elements up to columns 0, 1, 2, 3, 4 (in other words, the lag-1 difference of  $\mathbf{j}$  gives the number of non-zero elements for each column):

- $\mathbf{i} = (2, 3, 1, 2, 1, 3, 4, 1, 4);$

---

<sup>1</sup>“Special” means that there is a certain pattern for zero and / or non-zero elements. Sparse matrices of special sparse structure include for example, block diagonal matrices and banded matrices. In fact, this thesis has already seen these sparse matrices. The  $\mathbf{S}_\lambda$  and its additive components  $\tilde{\mathbf{S}}_l$  in REML estimation (see equation (4.3) in §4.3) are block diagonals, and the diagonal or tri-diagonal weight matrix  $\mathbf{W}$  are banded matrices.



- $\mathbf{j} = (0, 2, 4, 7, 9)$ ;
- $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9)$ .

$\mathbf{j}$  is often called a vector of *column pointers*, because the first non-zero element (and its row index) in column  $v$  can be located by skipping the first  $\mathbf{j}(v)$  values in  $\mathbf{x}$  (and  $\mathbf{i}$ ).

Essentially the triplet storage format still interprets the example  $\bar{\mathbf{W}}$  as a  $4 \times 4$  matrix, although only making reference to its non-zero elements. This makes a sparse matrix human-readable, but is not convenient for practical operations with the sparse matrix like matrix-vector and matrix-matrix multiplications. These operations access the sparse matrix by row or column, thus require an easy way to extract non-zero values in  $\mathbf{x}$  into row or column vectors. The CCR and CCS format are just for this purpose. For example, it is straightforward to extract column vectors (as compressed sparse vectors) from CCS format of the example  $\bar{\mathbf{W}}$ :

- skipping  $\mathbf{j}(1) = 0$  elements in  $\mathbf{i}$  and  $\mathbf{x}$  and reading in  $\mathbf{j}(2) - \mathbf{j}(1) = 2$  elements extracts column 1 as a compressed sparse vector:  $\mathbf{i}^{[1]} = (2, 3)$ ,  $\mathbf{x}^{[1]} = (x_1, x_2)$ ;
- skipping  $\mathbf{j}(2) = 2$  elements in  $\mathbf{i}$  and  $\mathbf{x}$  and reading in  $\mathbf{j}(3) - \mathbf{j}(2) = 2$  elements extracts column 2 as a compressed sparse vector:  $\mathbf{i}^{[2]} = (1, 2)$ ,  $\mathbf{x}^{[2]} = (x_3, x_4)$ ;
- skipping  $\mathbf{j}(3) = 4$  elements in  $\mathbf{i}$  and  $\mathbf{x}$  and reading in  $\mathbf{j}(4) - \mathbf{j}(3) = 3$  elements extracts column 3 as a compressed sparse vector:  $\mathbf{i}^{[3]} = (1, 3, 4)$ ,  $\mathbf{x}^{[3]} = (x_5, x_6, x_7)$ ;
- skipping  $\mathbf{j}(4) = 7$  elements in  $\mathbf{i}$  and  $\mathbf{x}$  and reading in  $\mathbf{j}(5) - \mathbf{j}(4) = 2$  elements extracts column 4 as a compressed sparse vector:  $\mathbf{i}^{[4]} = (1, 4)$ ,  $\mathbf{x}^{[4]} = (x_8, x_9)$ .

In **bamboos**, the primary operation with a sparse matrix is matrix-matrix multiplication like  $\bar{\mathbf{A}}'\bar{\mathbf{W}}$  and  $\bar{\mathbf{W}}\bar{\mathbf{B}}$ . Using CCS format for  $\bar{\mathbf{W}}$  is thus more helpful. In the next section, I will sketch algorithms for both matrix multiplications.

### Matrix multiplication with sparse $\bar{\mathbf{W}}$

Matrix-matrix multiplications  $\tilde{\mathbf{A}} = \bar{\mathbf{A}}'\bar{\mathbf{W}}$  and  $\tilde{\mathbf{B}} = \bar{\mathbf{W}}\bar{\mathbf{B}}$  in **bamboos** are not typical sparse matrix multiplications, because of the following.

- Only  $\bar{\mathbf{W}}$  is sparse, while  $\bar{\mathbf{A}}$  or  $\bar{\mathbf{B}}$  is generally dense. There are exceptions of course. The packed design matrix for an i.i.d. random effect is an identity matrix, which is sparse. However, in this situation there is no need to even perform matrix multiplication. In practice this special case can be recognized and dealt with separately.
- The resulting matrix  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{B}}$  are generally dense. This comes from the fact that each row or column of  $\bar{\mathbf{W}}$  has at least one non-zero elements (To see this, just note that the discretization index vectors  $\mathbf{k}_{\bar{\mathbf{A}}}$  and  $\mathbf{k}_{\bar{\mathbf{B}}}$  contain all positive integers from 1 to  $m_{\bar{\mathbf{A}}}$  and from 1 to  $m_{\bar{\mathbf{B}}}$ ). Then in  $\tilde{\mathbf{B}} = \bar{\mathbf{W}}\bar{\mathbf{B}}$ , each row of  $\tilde{\mathbf{B}}$  is a linear combination of rows of  $\mathbf{B}$ . Due to the previous fact, for any row of  $\tilde{\mathbf{B}}$ , there will be at least one row of  $\mathbf{B}$  added to that row. Thus if  $\mathbf{B}$  is dense,  $\tilde{\mathbf{B}}$  is dense. That  $\tilde{\mathbf{A}}$  is dense follows similarly.

These two characteristics make both the matrix-matrix multiplications very easy to program (without using any standard sparse linear algebra libraries). For example, consider programming  $\tilde{\mathbf{A}} = \bar{\mathbf{A}}'\bar{\mathbf{W}}$ , or  $\tilde{\mathbf{A}}(u, v) = \sum_k \bar{\mathbf{A}}(k, u)\bar{\mathbf{W}}(k, v)$ , then

<pre> initialize the resulting matrix <math>\tilde{\mathbf{A}} = \text{zeros}(p_{\tilde{\mathbf{A}}}, m_{\tilde{\mathbf{B}}})</math> for <math>v = 1 : m_{\tilde{\mathbf{B}}}</math>     number of non-zero elements in <math>\bar{\mathbf{W}}(:, v)</math>     <math>r = \mathbf{j}(v+1) - \mathbf{j}(v)</math>     offset for reading <math>\mathbf{i}</math> and <math>\mathbf{x}</math>     <math>s = \mathbf{j}(v)</math>     for <math>u = 1 : p_{\tilde{\mathbf{A}}}</math>         initialize an accumulator         <math>a = 0</math>         loop over non-zero entries of <math>\bar{\mathbf{W}}(:, v)</math>         for <math>k = 1 : r</math>             <math>a += \bar{\mathbf{A}}(\mathbf{i}(s+k), u)\mathbf{x}(s+k)</math>         write-back         <math>\tilde{\mathbf{A}}(u, v) = a</math> </pre>	<pre> initialize the resulting matrix <math>\tilde{\mathbf{B}} = \text{zeros}(m_{\tilde{\mathbf{A}}}, p_{\tilde{\mathbf{B}}})</math> for <math>v = 1 : p_{\tilde{\mathbf{B}}}</math>     for <math>k = 1 : m_{\tilde{\mathbf{B}}}</math>         linear combination coefficient         <math>b = \bar{\mathbf{B}}(k, v)</math>         number of non-zero elements in <math>\bar{\mathbf{W}}(:, k)</math>         <math>r = \mathbf{j}(k+1) - \mathbf{j}(k)</math>         offset for reading <math>\mathbf{i}</math> and <math>\mathbf{x}</math>         <math>s = \mathbf{j}(k)</math>         loop over non-zero entries of <math>\bar{\mathbf{W}}(:, k)</math>         for <math>u = 1 : r</math>             <math>\tilde{\mathbf{B}}(\mathbf{i}(s+u), v) += \mathbf{x}(s+u)b</math> </pre>
--	---

**Figure 9.9:** Matrix multiplication  $\tilde{\mathbf{A}} = \bar{\mathbf{A}}'\bar{\mathbf{W}}$  (left panel) and  $\tilde{\mathbf{B}} = \bar{\mathbf{W}}\bar{\mathbf{B}}$  (right panel) when  $\bar{\mathbf{W}}$  is sparse and stored with CCS format:  $\mathbf{i}$  is the vector of row index,  $\mathbf{j}$  is the vector of column pointers and  $\mathbf{x}$  is the vector of  $\bar{\mathbf{W}}$ 's non-zero entries.

- it is sufficient to only check if  $\bar{\mathbf{W}}(k, v)$  is zero or not when computing the inner product, because for any non-zero  $\bar{\mathbf{W}}(k, v)$ , the corresponding  $\bar{\mathbf{A}}(k, u)$  is always non-zero;
- there is no need to trace the computation to judge whether  $\tilde{\mathbf{A}}$  is dense or sparse hence what storage format should be used for it, because it is known to be dense.

Figure 9.9 sketches the algorithms for both matrix multiplications. Note that  $\tilde{\mathbf{B}} = \bar{\mathbf{W}}\bar{\mathbf{B}}$  or  $\tilde{\mathbf{B}}(u, v) = \sum_k \bar{\mathbf{W}}(u, k)\bar{\mathbf{B}}(k, v)$  is not programmed with  $v$ - $u$ - $k$  (outermost to innermost) loop ordering, but  $v$ - $k$ - $u$  loop ordering. This is because the former arrangement requires CRS format not CCS format.

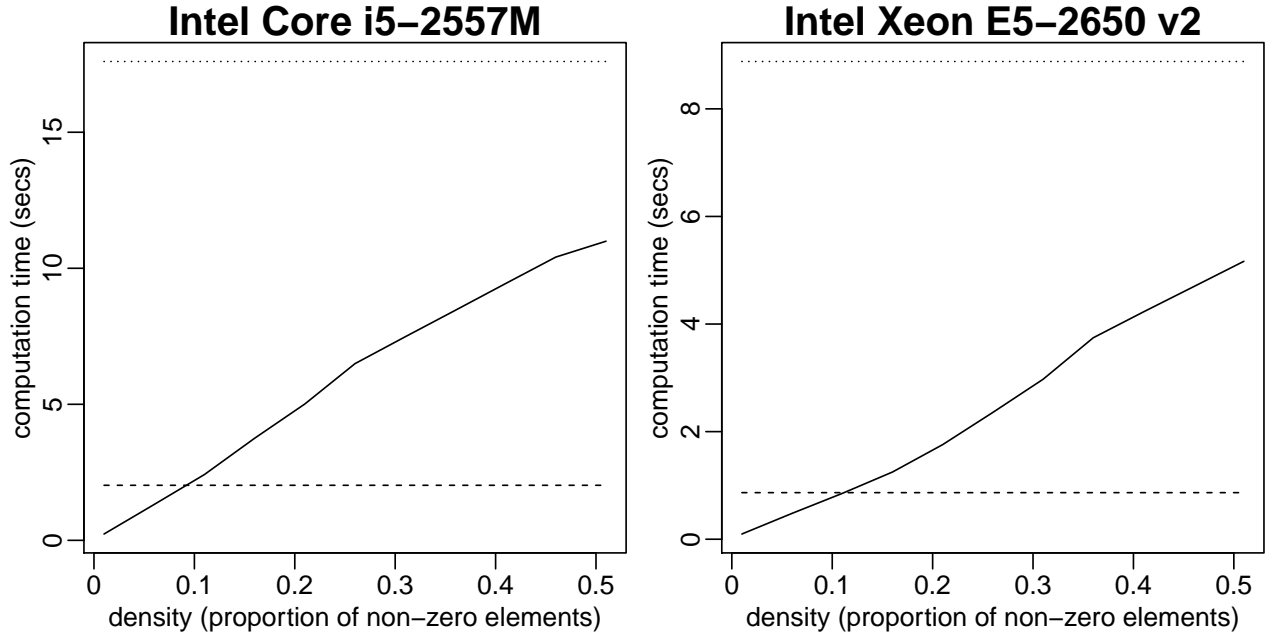
### “Dense-sparse switch” in bamboos

The nested sorting, run-length encoding and bin aggregation yields CCS format for  $\bar{\mathbf{W}}$ , making it equally straightforward to construct  $\bar{\mathbf{W}}$  as dense or sparse, hence it is possible to leave the construction option as a run-time decision. The critical quantity for the decision is  $K_{st}^{[1]}$ , the number of non-zero elements in  $\bar{\mathbf{W}}$  (the superscript is a reminder that I have been discussing diagonal  $\mathbf{W}$ ), which is known from the last element of the vector of column pointers. A threshold value between  $[0, 1]$  can be set for  $\mathcal{D}_{st}^{[1]} = K_{st}^{[1]}/m_{\tilde{\mathbf{A}}}m_{\tilde{\mathbf{B}}}$ , the proportion of non-zero elements in  $\bar{\mathbf{W}}$ , which is a measure of the density of the matrix.

- If  $\mathcal{D}_{st}^{[1]}$  is above this threshold,  $\bar{\mathbf{W}}$  should be constructed as dense (see Figure 9.7). Subsequent dense matrix-matrix multiplications can use level-3 BLAS for high-performance computing.
- If  $\mathcal{D}_{st}^{[1]}$  drops below this threshold,  $\bar{\mathbf{W}}$  should be constructed as sparse. A length- $K_{st}^{[1]}$  vector is initialized, then bin aggregation will aggregate weights  $\dot{\mathbf{w}}_{gh}$  onto this vector. A sparse matrix multiplication is first computed using algorithms in Figure 9.9, then all subsequent matrix-matrix multiplications are dense and level-3 BLAS can be used for high-performance computing.

The threshold therefore acts like a switch between dense construction and sparse construction. Fixing the threshold at 0 and 1 respectively enforce dense construction and sparse construction.

The threshold is a performance-related factor. It is not difficult to understand the following cases.



**Figure 9.10:** An empirical search for a reasonable threshold on Intel Core i5-2557M and Intel Xeon E5-2650 v2 (see Appendix A for hardware information).  $\bar{\mathbf{W}}$  is constructed as a  $2000 \times 2000$  random sparse matrix, with density  $\mathcal{D}_{st}^{[1]}$  (x-axis) taking trial values 0.01, 0.06, 0.11, ..., 0.51. In both panels, the y-axis is the mean computation time of  $\mathbf{W}\bar{\mathbf{B}}$  and  $\bar{\mathbf{A}}'\mathbf{W}$ , where  $\bar{\mathbf{A}}$  and  $\bar{\mathbf{B}}$  are  $2000 \times 2000$  random dense matrices. Three computation methods are used: solid line for sparse method, i.e., algorithms in Figure 9.9; dashed line for dense method with OpenBLAS; dotted line for dense method with reference BLAS. On both machines, the solid line intersects the dashed line near 0.1, thus this is a reasonable threshold value for dense-sparse switch when an optimized BLAS is used.

- When  $\mathcal{D}_{st}^{[1]}$  is big, i.e., when  $\bar{\mathbf{W}}$  is realistically dense, doing matrix multiplication with the sparse method will induce great overhead due to index loopup, hence the dense method will outperform the sparse method.
- When  $\mathcal{D}_{st}^{[1]}$  is small, i.e., when  $\bar{\mathbf{W}}$  is realistically sparse, doing matrix multiplication with the dense method will waste substantial amount of time doing computation with zeros, hence the sparse method will outperform the dense method.

There must exist a crossover  $\mathcal{D}_{st}^{[1]}$ , and in theory this will be the optimal threshold. However, there is no other ways than practical benchmarking to make a proper guess on this value. In addition, this value is likely to depend on the following factors:

1. the machine where computation is performed;
2. the BLAS distribution that is used (it is not difficult to predict that with an optimized BLAS in place of the reference BLAS, it is more difficult for the sparse method to beat the dense method).

Figure 9.10 conducts an empirical search for a reasonable threshold on Intel Core i5-2557M and Intel Xeon E5-2650 v2 (see Appendix A for hardware information). It appears that 0.1 is a reasonable choice when an optimized BLAS (OpenBLAS in the experiment) is used. This threshold is currently adopted in **bamboos**.

### 9.3.4 Tri-diagonal $\mathbf{W}$

Methods in the last section (§9.3) still works when the weight matrix  $\mathbf{W}$  is tri-diagonal weight, by dealing its three diagonals separately. For example, the nested sorting will be applied to  $(\mathbf{k}_{\bar{\mathbf{A}}}, \mathbf{k}_{\bar{\mathbf{B}}}, \mathbf{w}_0)$ ,

$(\mathbf{k}_A(1 : (n-1)), \mathbf{k}_B(2 : n), \mathbf{w}_1)$  and  $(\mathbf{k}_A(2 : n), \mathbf{k}_B(1 : (n-1)), \mathbf{w}_{-1})$  independently. Bin aggregation is applied in each case, producing three contingency matrices. They are then added together for the final  $\bar{\mathbf{W}}$ .

## 9.4 Experimenting bamboos on daily logBS model

Previous sections have been a thorough demonstration of the design of **bamboos**. Now it is time to assess how worthwhile **bamboos** is in practical large GAM fitting. In this section, I will experiment **bamboos** on the test daily logBS model, model 7.1, and make performance comparison with the old discrete algorithm.

Recall from §8.6 that to compute  $\mathbf{X}'\mathbf{W}\mathbf{X}$ ,

- 16-threaded “online” pseudo QR reduction with OpenBLAS takes 50 minutes;
- 16-threaded old discrete computation takes 34 minutes.

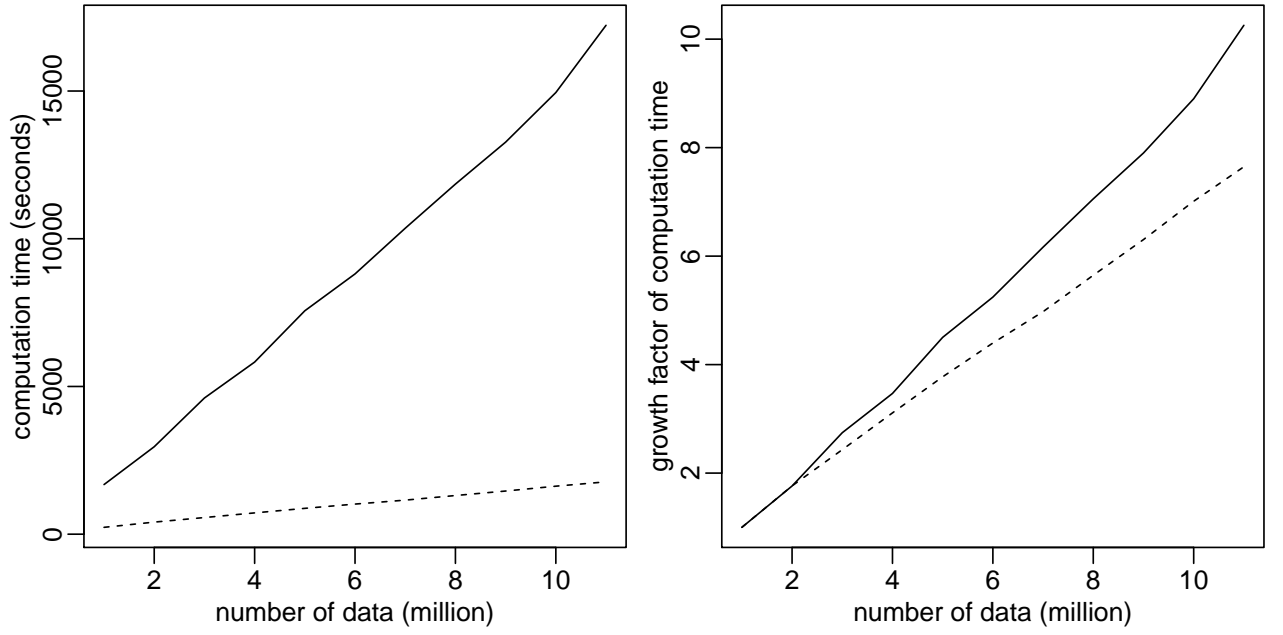
Application of **bamboos** shows that it takes 12 minutes, with just a single thread. That is, **bamboos** is so fast that even its serial computing performance has beaten previous methods with the aid of massively parallel computing. **bamboos** algorithms can be further parallelized, by parallelizing the computation within  $T \times S$ ,  $S \times T$  and  $T \times T$  blocks because of the following. Referring to Figure 9.4, different iterations are embarrassingly parallel; in addition, they have the same computational complexity so load balancing is guaranteed. At it turns out, **bamboos** has very good parallel scaling. Using 4 threads, computation of  $\mathbf{X}'\mathbf{W}\mathbf{X}$  completes in 3.3 minutes. And there is no need to use more threads.

A big advantage of **bamboos**, is that the more data there are, the faster it is compared with previous computation methods. This is what has been theoretically concluded from the  $n \rightarrow \infty$  case analysis in §9.2.3. To verify this in practice, I will subsample (with replacement) the Monday logBS dataset for  $1 \times 10^6$ ,  $2 \times 10^6$ , ...,  $11 \times 10^6$  data and fit the test model 7.1. I will just compare **bamboos** (serial computing) with the old discrete method (16-threaded computing). The result is sketched in Figure 9.11. It can be seen that when data are 10 times as more, computation time for old discrete method is about 10 times as longer, but computation time for **bamboos** is just about less than 8 times as longer. As another test of **bamboos**, consider dropping the three-way interactions (three-margin tensor product splines)  $f_{16}$  and  $f_{17}$  from the test model, and redo the above experiment. This time it can be seen that when data are 15 times as more, computation time for **bamboos** is less than 7 times as longer.

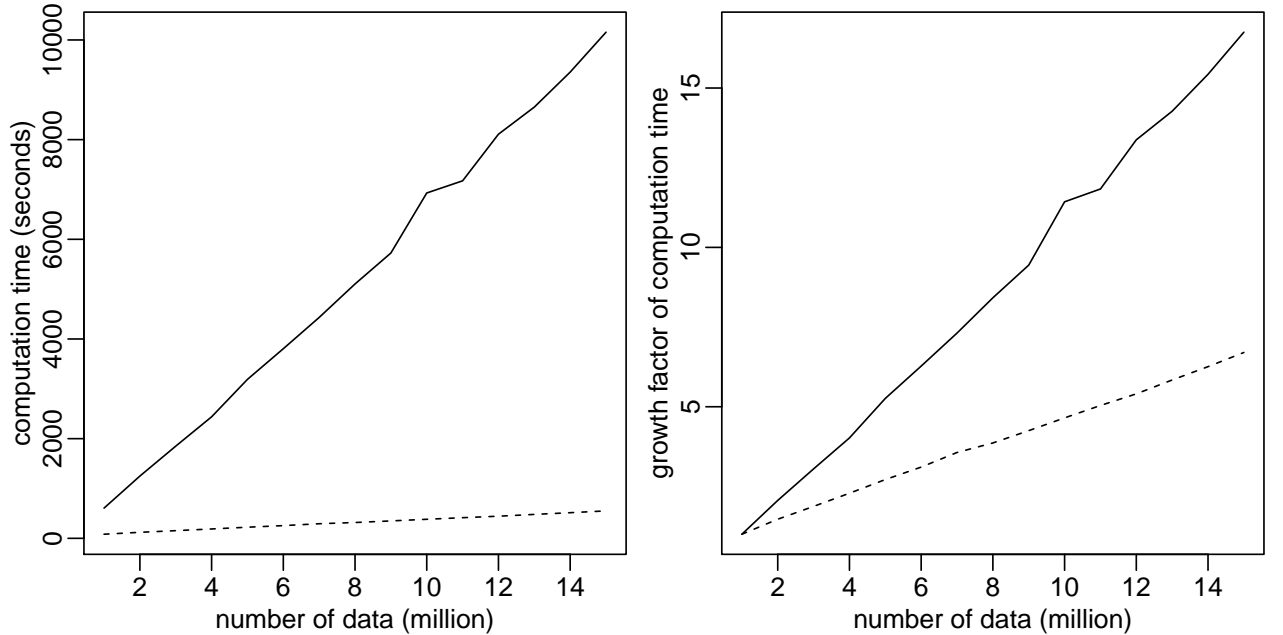
## 9.5 Summary

In this Chapter I have developed a new computational engine, **bamboos**, for the matrix cross-product  $\mathbf{X}'\mathbf{W}\mathbf{X}$ , when the submatrices of  $\mathbf{X}$  corresponding to individual GAM terms in a linear predictor  $\boldsymbol{\eta} = \mathbf{X}_0\boldsymbol{\beta}_0 + \sum_{i=1}^l \mathbf{X}_i\boldsymbol{\beta}_i$  are stored in a packed storage format from covariate discretization (described in §8.2). The method successfully boosts the previously successful discrete algorithms in Chapter 8 (and published in Wood et al. (2017)) by another magnitude.

**bamboos** is fast for the following reasons.



**Figure 9.11:** Benchmarking **bamboos** (serial computing) and old discrete method (16-threaded) on Intel E5-2650 v2 (see Appendix A for hardware information) by fitting test daily logBS model to  $1 \times 10^6$ ,  $2 \times 10^6$ , ...,  $11 \times 10^6$  data. The left panel shows the computation time for both methods (solid: old discrete; dashed: **bamboos**). For 1 million data, old discrete method takes 1680 seconds and **bamboos** takes 230 seconds; for 11 million data, old discrete method takes 17225 seconds and **bamboos** takes 1772 seconds. The right panel shows the growth factor of computation time for both methods (solid: old discrete; dashed: **bamboos**). It can be seen that when data are 10 times as more, computation time for old discrete method is about 10 times as longer, but computation time for **bamboos** is just about less than 8 times as longer.



**Figure 9.12:** Benchmarking **bamboos** (serial computing) and old discrete method (16-threaded) on Intel E5-2650 v2 (see Appendix A for hardware information) by fitting test daily logBS model (dropping  $f_{16}$  and  $f_{17}$ ) to  $1 \times 10^6$ ,  $2 \times 10^6$ , ...,  $15 \times 10^6$  data. The left panel shows the computation time for both methods (solid: old discrete; dashed: **bamboos**). For 1 million data, old discrete method takes 606 seconds and **bamboos** takes 82 seconds; for 15 million data, old discrete method takes 10155 seconds and **bamboos** takes 550 seconds. The right panel shows the growth factor of computation time for both methods (solid: old discrete; dashed: **bamboos**). It can be seen that when data are 15 times as more, computation time for old discrete method is about 15 times as longer, but computation time for **bamboos** is just about less than 7 times as longer.

1. Its fundamental design idea is superior to that of the old discrete methods. As is demonstrated in §9.1, **bamboos** algorithms are block-oriented, rich in matrix-matrix multiplications. By contrast, the old discrete algorithms are vector-oriented, rich in matrix-vector operations (like unpacking and aggregation). Therefore, **bamboos** has much better opportunity to benefit from optimized BLAS.
2. While the advantage of **bamboos** is not immediately clear from the upper bound of its FLOP count when compared with the old discrete algorithms (see §9.2), it gains substantial practical computational savings due to sparsity. Therefore, the speedup of **bamboos** is very often, huge, as is seen from the benchmarking in §9.4. Note that **bamboos** features a run-time decision between doing dense and sparse computations. This is very appealing, since users do not need to make judgement themselves. Other R packages like **glm4** that is successful in fitting large GLMs does require an exclusive user-specification on dense or sparse methods.
3. **bamboos** does not completely eliminate vector-oriented operations. The weights aggregation is a typical example, and in fact in the  $n \rightarrow \infty$  limiting case, it will dominate **bamboos** computations. However, the nested sorting, run-length encoding and bin aggregation strategy makes such operation as fast as possible by optimizing memory access. The practical gains of these optimizations can be seen by comparing **bamboos** with its adapted version implemented in **mgcv** from version 1.8-25 (2018-10-26): **bamboos** is 6 times as fast in computing  $\mathbf{X}'\mathbf{W}\mathbf{X}$  for the test daily logBS model, model 7.1.

Basically, **mgcv** has chosen a balance point between algorithmic efficiency and coding complexity. Admittedly, the nested sorting and sparse representation of  $\bar{\mathbf{W}}$  add substantial difficulty in coding. That **bamboos** tries to accurately determine the degree of sparsity / density of  $\bar{\mathbf{W}}$  (see quantities  $K_{st}^{[1]}$  and  $\mathcal{D}_{st}^{[1]}$ ), choose a near optimal dense-sparse switch, and optimize every memory access it sees is a heavy burden for its development. This is particularly true if  $\mathbf{W}$  is tri-diagonal. **mgcv** uses a simpler heuristic on the sparsity of  $\bar{\mathbf{W}}$  from the crude relationship between  $n$  and  $m_{\bar{A}_s}m_{\bar{B}_t}$ .

- If  $m_{\bar{A}_s}m_{\bar{B}_t} < n$ , the matrix is seen dense, and algorithm in the right panel of Figure 9.1a is used.
- If  $m_{\bar{A}_s}m_{\bar{B}_t} > n$ , the matrix is seen sparse. Yet **mgcv** opts to skip the explicit sparse construction of  $\bar{\mathbf{W}}$ , computing  $\bar{\mathbf{A}}'\bar{\mathbf{W}}$  or  $\bar{\mathbf{W}}\bar{\mathbf{B}}$  directly by matrix row aggregation.

In this way, the coding complexity is substantially reduced, but neither option above is optimal. For example, in the first case, the random memory access in weights aggregation is not be resolved, so the performance penalty observed in Figure 9.8 will have negative effect. In addition,  $m_{\bar{A}_s}m_{\bar{B}_t} < n$  is a too optimistic bet on that  $\bar{\mathbf{W}}$  is dense. For example, when  $m_{\bar{A}_s}m_{\bar{B}_t} = 0.95n$ ,  $\bar{\mathbf{W}}$  is realistically likely to be sparse hence using dense computation method will lose some efficiency. Nevertheless, the adapted implementation in **mgcv** preserves many good features of original **bamboos** design, thus is still significantly faster than the old discrete algorithms (hence very useful). A paper “Faster model matrix crossproducts for large generalized linear models with discretized covariates” has been submitted to *Statistics and Computing* based on this implementation.

It is worth pointing out that while **bamboos** is originally designed for fitting large GAMs, broadly speaking, it offers advantages for any regression model in which each covariate results in a term with several associate model matrix columns: models containing several factor variables are an obvious example beyond smooth regression models. However, there is still a long way to go if **bamboos** is to be developed into an R package (like **Matrix** package) that provides generic computation methods. At the moment, I am hoding an unpublished (not even on GitHub, etc) **bamboos** package for methods described in this Chapter.

# Discussion

What is this thesis about? What is my PhD about? It is very difficult to summarize. Pessimistically speaking, it could be that I have attempted many things, but only half of them were made into this thesis, among which I have only been truly successful with one thing: **bamboos**.

**bamboos** offers an extremely efficient computation method for large GAM (or GLM) fitting. Strictly speaking, **bamboos** is only a far end of GAM fitting method for large datasets. It is built upon at least two critical previous research work.

- The idea of breaking up GAM fitting into model matrix reduction and smoothing parameter selection, and the application of “online” algorithm to the former.
- The discovery of the discrete characteristic in Black Smoke dataset, the proposal to thus store data in a compressed format, and the initial set of discrete algorithms that are able to do simple operations like vector unpacking and vector aggregation on data stored in such format.

The first idea had just been around when I started my PhD. The **bam** function implemented it. In today’s view, it was a very crude implementation, with many computations arranged at R-level rather than C-level. But this is not surprising. Back at that stage, ensuring correctness of computations and validity of the estimation idea was the top priority. There were thus many opportunities to optimize it. It was only known that it would be successful to overcome the memory bottleneck imposed by large datasets and large model matrices, but how efficient the method was in large GAM fitting was unknown. In fact, many other aspects had not been discovered as well, like how Cholesky factorization may handle rank-deficiency, and how psuedo QR reduction performs compared with QR reduction.

The Black Smoke dataset served a perfect example for testing the effectiveness of this fitting method. It was soon discovered that it was even difficult to fit a model to annual logBS of 20000 - 30000 data and 3000 - 5000 parameters. In reality, a predictive model for annual Smoke concentration is of little practical use; it is hard to find correlation between annual mean pollution level and epidemic. Therefore, there was also a need to model logBS at a finer temporal resolution, which inevitably ended up with GAM fitting for substantially larger datasets. Meanwhile, some slowish preliminary model also suggested that there there exists complicated interactions between variables and a reasonable model would have to involve a high number of parameters.

Building a satisfying statistical model is one question; being able to estimate one is another, and in some sense, more important. Without being able to fit a model first, there is no way to assess the quality that model is and how it may be improved. Based on this acknowledgment, the focus of the research increasingly shifted from building Black Smoke models to developing fast large GAM fitting methods, and it went further and further on this route. Eventually, this thesis became very strong on the computational aspect, but very weak in statistical modelling aspect.

Chapters of this thesis have been organized to highlight the early use of optimized BLAS for initial performance improvement of GAM fitting. However, the true story is, an optimized BLAS had not

been experimented even after the discrete algorithms were proposed, implemented and published. I started formal investigation on optimized BLAS when I tried to improve the performance of pivoted Cholesky factorization. Although after 6 - 12 months' investigation I was not able to roll out an implementation that consistently outperforms existing computation methods, I have learned a lot about how an optimized BLAS is designed, how to write fast C code, and how to read assembly output by a compiler and even write some assembly code for speed. Only by this stage, I then felt more competent in analyzing the performance of `mgcv`. What has been presented in Chapter 5 is really the result of my revisit to `bam`. I later found that the discrete method in Chapter 8, while substantially faster than “online” algorithms implemented with reference BLAS, is not competitive with that implemented with OpenBLAS. This greatly raised my curiosity, and I started to investigate discrete algorithms in depth again.

Note that it was already in the 5-th year of my PhD when I revisited the discrete algorithms. So while I was supposed to focus on writing my thesis, telling about the story of how the discrete algorithm was found and how daily Black Smoke model was built, I was spending substantial amount of time in researching the discrete algorithm. Initially I applied caching to a few steps of these discrete algorithms, but the gains was not impressive. I understood why: the old algorithms was vector-oriented, hard to attain high performance. The deadline for submitting my thesis was getting close, then I have to stop further investigation.

Personally speaking, I did not find anything I could be proud of when writing my thesis. The original idea of discrete algorithm was not proposed by me; my Black Smoke models were also very weak. Therefore, I did not have confidence that I could obtain a degree. The thesis was written very short, primarily because I did not see anything worthwhile to write (at least in my own opinion). To make my thesis longer, I tried hard, adding in many immature materials like using knots placement to suppress spatial-temporal prediction bias, developing some toy models for the dropout process of stations. However, these are not what I really want to do.

A month before the submission deadline, I suddenly got some idea on how to arrange discrete algorithms in block-oriented fashion. I was very excited. I was sure this is going to be a novel idea. The basic idea of these algorithms is not very complicated; however, the nested sorting, run-length encoding, bin-aggregation, handling sparsity turns out complicated. But the more complicated it is, the more worthwhile it is to do. Before submission of the thesis I was able to establish theoretical superiority of these methods. The first full implementation was not ready until my viva was near. I decided to place the development of these fast algorithms as my prominent contributions over the past 7 years.

I have been too narrowly focused on computational issues, and the building statistical models have been overlooked to some extent. In the revision / correction stage of the thesis, more prudent model development for logBS were added. These models turned out much more complicated than those in the original thesis. I did learn plenty of things in this process, for example,

- I implemented the golden-section search so that GAMs can be estimated together with AR(1) autocorrelation.
- I also learned to thin data to alleviate or even eliminate the effect of temporal autocorrelation in building GAMs for exploratory analysis.
- I also learned plenty of visualization methods for model checking.

At the same time, I also realized that it is very difficult to adequately model spatial autocorrelation in the data. I concluded that splines is more useful to remove long-range spatial trend, but less competent in modelling short range variability.

These revised model development altered my opinion of building daily logBS model. It is even difficult



to build an adequate model for thinned data (weekly logBS from each day of week) without temporal correlation, there is hardly any sense to move on further to a full daily model. If I would ever revisit Black Smoke models again, I would highly recommend first getting more daily covariate variables. The exploratory analysis methods presented in Chapter 3 will still be valid.

I honestly admitted that I am incompetent in building statistical models. However, I have provided fast model fitting methods, offering opportunities for those competent to investigate the dataset. I would still devote my future work on computational issues. Compared with model building, I feel that these matters are more likely to be firmly answered. Life is too short; I want to do something that is 100% answerable.

# Appendix A

## Hardware information

Although the computational methods reviewed and to be developed in this thesis are not targeted on any specific machine, there is still a need to introduce machines that are used for computations, profiling and benchmarking. Three machines are used in this thesis, and all of them come with Linux OS.

- Intel Core i5-2557M Sandy Bridge laptop. It has two CPU cores, each of which can work at a sustainable frequency of 1.1GHz. 4GB RAM is available, with 32KB L1 data cache for each core. There are two memory channels in total, each supporting DDR3 DIMM at 1333MHz, giving a maximum FSB bandwidth of 10.664 GB/s per channel. OpenBLAS is available on this machine for matrix computations;
- Intel Xeon E5-2650 v2 Ivy Bridge workstation. It has 16 CPU cores, each of which can work at a sustainable frequency of 2.6GHz. 128GB RAM is available, with 32KB L1 data Cache for each core. There are four memory channels in total, each supporting DDR3 DIMM at 1866 MHz, giving a maximum FSB bandwidth of 14.928 GB/s per channel. This machine is also the only one with NUMA feature, as its 16 cores sit on two CPU sockets. This machine is a node of University of Bath's HPC cluster *Balena*, launched in 2015. Both Intel MKL and OpenBLAS are available on this machine for matrix computations;
- Intel Core M-5Y71 Broadwell laptop. It has two CPU cores, each of which can work at a sustainable frequency of 1.3GHz. 8GB RAM is available, with 32KB L1 data cache for each core. There are two memory channels in total, each supporting LPDDR3 DIMM at 1600 MHz, giving a maximum FSB bandwidth of 12.8 GB/s per channel. OpenBLAS is available on this machine for matrix computations.

“Sustainable frequency” is a frequency that a CPU can sustain at during long-time intensive computations without overheating, in both single-threading and multi-threading application. All these machines come with Intel’s “turbo boost” technology, which enables a higher CPU frequency for a short burst of time in single-threading application. Core i5-2557M has a “turbo” range 1.7 ~ 2.7GHz, Xeon E5-2650 v2 has a “turbo” range 2.6 ~ 3.4GHz, Core M-5Y71 has a “turbo” range 1.3 ~ 2.9 GHz. CPU frequency in “turbo” range is not sustainable, as the excessive amount of heat it generates would eventually cause dynamic CPU frequency scaling to reduce the frequency. “Turbo boost” must be disabled to rule out the impact of unstable CPU frequency in benchmarking. Note that on Core i5-2557M, simply disabling “turbo boost” is not enough. The maximum 1.7GHz frequency below “turbo” range can still cause CPU temperature to reach the critical 86°C within 1 minutes. It turns out that 1.1GHz is its sustainable frequency.

Originally all machines have the capability for scheduling hyper-threading, but it is later disabled. This technology allows two threads to be scheduled on each CPU core, so it appears to the operating

system as well as users that there are two logical CPU cores per physical CPU core. However, logical cores have CPU affinity of 1, sharing all physical resources like CPU registers and CPU cache. In computationally intensive tasks like dense matrix computations, hyper-threading would hammer practical parallel scalability, because logical cores would be busy competing for physical resources and having two of them is hardly any better than having just one. Hyper-threading is more promising in sparse matrix computations which is more memory intensive, as if one thread has a cache miss in reading data from RAM, the other thread may have a cache hit so the CPU core can avoid being idle.

# Bibliography

- Anderson, E., Bai, Z., Bischof, C., Blackford, L. S., Demmel, J., Dongarra, J. J., Du Croz, J., Hammarling, S., Greenbaum, A., McKenney, A., and Sorensen, D. (1999). *LAPACK Users' Guide*. Society for Industrial and Applied Mathematics, 3 edition.
- Anderson, E., Bai, Z., Bischof, C., Demmel, J. W., Dongarra, J. J., Croz, J. D., Greenbaum, A., Hammarling, S., McKenney, A., and Sorensen, D. C. (1990). LAPACK: A portable linear algebra library for high-performance computers. Technical Report 20, LAPACK Working Note.
- Anderson, E. and Dongarra, J. J. (1990). Implementation guide for LAPACK. Technical Report 18, LAPACK Working Note.
- Augustin, N. H., Musio, M., von Wilpert, K., Kublin, E., Wood, S. N., and Schumacher, M. (2009). Modeling spatiotemporal forest health monitoring data. *Journal of the American Statistical Association*, 104(487):899–911.
- Belitz, C. and Lang, S. (2008). Simultaneous selection of variables and smoothing parameters in structured additive regression models. *Computational Statistics and Data Analysis*, 53(1):61–81.
- Bischof, C. and Loan, C. V. (1987). The wy representation for products of householder matrices. *SIAM Journal on Scientific and Statistical Computing*, 8(1):s2–s13.
- Bowman, A. and Azzalini, A. (1997). *Applied Smoothing Techniques for Data Analysis: The Kernel Approach with S-Plus Illustrations*. Oxford Statistical Science Series. OUP Oxford.
- Breslow, N. E. and Clayton, D. G. (1993). Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association*, 88(421):9–25.
- Burylov, I., Chuvelev, M., et al. (2007). Intel performance libraries: multi-core-ready software for numeric-intensive computation. *Intel technology journal*, 11(4):299–308.
- Ciocco, A. and Thompson, D. (1961). A follow-up of donora ten years after: methodology and findings. *Am J Public Health Nations Health*, 51:155–164.
- Claeskens, G., Krivobokova, T., and Opsomer, J. D. (2009). Asymptotic properties of penalized spline estimators. *Biometrika*, 96(3):529–544.
- Cleveland, W. S. (1979). Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association*, 74(368):829–836.
- Cleveland, W. S. and Devlin, S. J. (1988). Locally weighted regression: An approach to regression analysis by local fitting. *Journal of the American Statistical Association*, 83(403):596–610.
- Clifton, M. (1964). Air Pollution. *Proceedings of the Royal Society of Medicine*, 57(7):615–618.
- Crainiceanu, C., Ruppert, D., and Wand, M. (2005). Bayesian analysis for penalized spline regression using winbugs. *Journal of Statistical Software, Articles*, 14(14):1–24.
- Craven, P. and Wahba, G. (1979). Smoothing noisy data with spline functions. *Numerische Mathematik*, 31:377–403.
- Croz, J. D. and Higham, N. J. (1992). Stability of methods for matrix inversion. *IMA Journal of Numerical Analysis*, 12(1):1–19.

- Dadvand, P., Rushton, S., et al. (2011). Using spatio-temporal modelling to predict long-term exposure to Black Smoke at fine spatial and temporal scale. *Atmosphere Environment*, 45(3):659–664.
- de Boor, C. (1978). *A Practical Guide to Splines*. Applied Mathematical Sciences. Springer-Verlag New York.
- Demmel, J. W., Dongarra, J. J., Croz, J. D., Greenbaum, A., Hammarling, S., and Sorensen, D. C. (1987). Prospectus for the development of a linear algebra library for high-performance computers. Technical Report 1, LAPACK Working Note.
- Demmel, J. W., Hoemmen, M., Hida, Y., and Riedy, E. J. (2009). Nonnegative diagonals and high performance on low-profile matrices from householder qr. *SIAM Journal on Scientific Computing*, 31(4):2832–2841.
- Dongarra, J. J., Croz, J. D., Hammarling, S., and Duff, I. (1990). A set of level 3 basic linear algebra subprograms. *ACM transactions on mathematical software*, 16(1):1–17.
- Dongarra, J. J., Croz, J. D., Hammarling, S., and Hanson, R. J. (1988). An extended set of basic of fortran basic linear algebra subprograms. *ACM transactions on mathematical software*, 14(1):1–17.
- Drepper, U. (2007). What every programmer should know about memory. Technical report, Red Hat, Inc. <https://people.freebsd.org/~lstewart/articles/cpumemory.pdf>.
- Duchon, J. (1977). *Constructive theory of functions of several variables*, volume 571, chapter Splines minimizing rotation-invariant semi-norms in Solobev spaces, pages 85–100. Springer Berlin Heidelberg.
- Duersch, J. A. and Gu, M. (2017). Randomized qr with column pivoting. *SIAM Journal on Scientific Computing*, 39(4):C263–C291.
- Eilers, P. H. C. and Marx, B. D. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science*, 11(2):89–121.
- Eilers, P. H. C. and Marx, B. D. (2003). Multivariate calibration with temperature interaction using two-dimensional penalized signal regression. *Chemometrics and intelligent laboratory systems*, 66(2):159–174.
- Enea, M. (2009). Fitting linear models and generalized linear models with large data sets in R. In *book of short papers, conference on "Statistical Methods for the analysis of large data-sets"*, pages 411–414. Italian Statistical Society, Chieti-Pescara,.
- Epperson, J. F. (1987). On the runge example. *American Mathematical Monthly*, 94(4):329–341.
- Fahrmeir, L. and Lang, S. (2001). Bayesian inference for generalized additive mixed models based on markov random field priors. *Applied Statistics*, 50:201–220.
- Fan, J. and Gijbels, I. (1996). *Local Polynomial Modelling and Its Applications*, volume 66 of *Mono-graphs on Statistics and Applied Probability*. Chapman and Hall/CRC.
- Fanshawe, T. R., Diggle, P. J., et al. (2008). Modelling spatio-temporal variation in exposure to particulate matter: a two-stage approach. *Environmetrics*, 19(6):549–566.
- Faraway, J. J. (2004). *Linear models with R*. Texts in Statistical Science. Chapman and Hall/CRC.
- Firket, J. (1936). Fog along the Meuse valley. *Transactions of the Faraday Society*, 32:1191–1194.
- Golub, G. H. and Loan, C. F. V. (2013). *Matrix Computations*. The Johns Hopkins University Press, Baltimore, 4 edition.
- Goto, K. and Geijn, R. V. D. (2008a). Anatomy of high-performance matrix multiplication. *ACM transactions on mathematical software*, 34(3).
- Goto, K. and Geijn, R. V. D. (2008b). High-performance implementation of the level-3 BLAS. *ACM transactions on mathematical software*, 35(1).

- Green, P. and Silverman, B. (1994). *Nonparametric Regression and Generalized Linear Models*, volume 58 of *Monographs on Statistics and Applied Probability*. Chapman and Hall.
- Gu, C. (1992). Cross-validating non-gaussian data. *Journal of Computational and Graphical Statistics*, 1(2):169–179.
- Gu, C. and Kim, Y. J. (2002). Penalized likelihood regression: general approximation and efficient approximation. *Canadian Journal of Statistics*, 34(4):619–628.
- Gulliver, J., Morris, C., et al. (2011). Land use regression modeling to estimate historic (1962-1991) concentrations of Black Smoke and Sulphur Dioxide for Great Britain. *Environmental Science and Technology*, 45:3526–3532.
- Gunnels, J. A., Henry, G. M., and Geijn, R. A. (2001). *ICCS 2001: International Conference San Francisco, CA, USA, May 28–30, 2001 Proceedings, Part I*, chapter A Family of High-Performance Matrix Multiplication Algorithms, pages 51–60. Springer Berlin Heidelberg.
- Hall, P. and Opsomer, J. D. (2005). Theory for penalised spline regression. *Biometrika*, 92(1):105–118.
- Handcock, M. S., Meier, K., and Nychka, D. (1994). Comment. *Journal of the American Statistical Association*, 89(426):401–403.
- Harville, D. A. (1997). *Matrix algebra from a statistician’s perspective*. Springer.
- Hastie, T. and Tibshirani, R. (1986). Generalized additive models (with discussion). *Statistical Science*, 1:297–318.
- Hastie, T. and Tibshirani, R. (1990). *Generalized Additive Models*. Chapman and Hall.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *Kernel Smoothing Methods*, pages 191–218. Springer New York, New York, NY.
- Helwig, N. E. and Ma, P. (2016). Smoothing spline ANOVA for super large samples: scalable computation via rounding parameters. *preprint arXiv:1602.05208*.
- Hennessy, J. L. and Patterson, D. (2011). *Computer Architecture: A Quantitative Approach*. The Morgan Kaufmann Series in Computer Architecture and Design. Elsevier, 5 edition.
- Kågström, B., Ling, P., and Loan, C. V. (1998). Gemm-based level 3 BLAS: High performance model implementations and performance evaluation benchmark. *ACM transactions on mathematical software*, 24(3):268–302.
- Kammann, E. E. and Wand, M. P. (2003). Geoadditive models. *Applied Statistics*, 52(1):1–18.
- Kauermann, G., Krivobokova, T., and Fahrmeir, L. (2009). Some asymptotic results on generalized penalized spline smoothing. *Journal of the Royal Statistical Society: Series B*, 71(2):487–503.
- Kimeldorf, G. S. and Wahba, G. (1970). A correspondence between bayesian estimation on stochastic processes and smoothing by splines. *The Annals of Mathematical Statistics*, 41(2):495–502.
- Kowarschik, M. and Weiß, C. (2003). *An Overview of Cache Optimization Techniques and Cache-Aware Numerical Algorithms*, pages 213–232. Springer Berlin Heidelberg, Berlin, Heidelberg.
- Kress, R. (1996). *Numerical Analysis*, volume 181 of *Graduate Texts in Mathematics*. Springer-Verlag New York.
- Krivobokova, T. and Kauermann, G. (2007). A note on penalized spline smoothing with correlated errors. *Journal of the American Statistical Association*, 102(480):1328–1337.
- Lam, M. D., Rothberg, E. E., and Wolf, M. E. (1991). The cache performance and optimizations of blocked algorithms. *SIGPLAN Notices*, 26(4):63–74.
- Lang, S., Umlauf, N., Wechselberger, P., et al. (2014). Multilevel structured additive regression. *Statistics and Computing*, 24(2):223–238.

- Lawson, C. L., Hanson, R. J., Kincaid, D. R., and Krogh, F. T. (1979). Basic linear algebra subprograms for fortran usage. *ACM transactions on mathematical software*, 5(3):308–323.
- Lee, D.-J. and Durbán, M. (2011). P-spline anova-type interaction models for spatio-temporal smoothing. *Statistical Modelling*, 11(1):49–69.
- Li, Y. and Ruppert, D. (2008). On the asymptotics of penalized splines. *Biometrika*, 95(2):415–436.
- Lindgren, F. and Rue, H. (2015). Bayesian spatial modelling with r-inla. *Journal of Statistical Software*, 63(19):1–25.
- Lindgren, F., Rue, H., and Lindström, J. (2011). An explicit link between gaussian fields and gaussian markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 73(4):423–498.
- Loader, A. (2002). *Instruction manual: UK Smoke and Sulphur Dioxide Network*. Netcen, AEA Technology, Culham Science Centre.
- Lucas, C. (2004). LAPACK-style codes for level 2 and 3 pivoted cholesky factorizations. Technical Report 161, LAPACK Working Note.
- Martinsson, P.-G., Ort, G. Q., Heavner, N., and van de Geijn, R. (2017). Householder qr factorization with randomization for column pivoting (hqrrp). *SIAM Journal on Scientific Computing*, 39(2):C96–C115.
- Marx, B. D. and Eilers, P. H. C. (1998). Direct generalized additive modeling with penalized likelihood. *Computational Statistics and Data Analysis*, 28:193–209.
- Miller, A. J. (1992). Algorithm AS 274: Least Squares Routines to Supplement Those of Gentleman. *Journal of the Royal Statistical Society: Series C*, 41(2):458–478.
- Ministry of Health (1954). *Mortality and morbidity during the London fog of December 1952*. HMSO, London.
- Nadaraya, E. A. (1964). On estimating regression. *Theory of Probability & Its Applications*, 9(1):141–142.
- Nocedal, J. and Wright, S. (2006). *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer.
- Opsomer, J., Wang, Y., and Yang, Y. (2001). Nonparametric regression with correlated errors. *Statistical Science*, 16(2):134–153.
- Parker, P. and Rice, J. (1985). Discussion of silverman (1985). *Journal of the Royal Statistical Society: series B*, 47(1):40–42.
- Patterson, H. D. and Thompson, R. (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika*, 58(3):545–554.
- Perry, M. and Hollis, D. (2006). *The generation of monthly gridded datasets for a range of climatic variables over the United Kingdom*. Met Office, FitzRoy Road, Exeter, Devon, EX1 3PB, United Kingdom.
- Perry, M., Hollis, D., and Elms, M. (2009). *The Generation of Daily Gridded Datasets of Temperature and Rainfall for the UK*. National Climate Information Centre, Met Office, FitzRoy Road, Exeter, Devon, EX1 3PB, United Kingdom.
- Plummer, M. (2003). JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. In *Proceedings of the 3rd International Workshop on Distributed Statistical Computing (DSC 2003)*.
- Quintana-Ortí, G., Sun, X., and Bischof, C. H. (1998). A BLAS-3 version of the QR factorization with column pivoting. *SIAM Journal on Scientific Computing*, 19(5):1486–1494.

- Reinsch, C. H. (1967). Smoothing by spline functions. *Numerische Mathematik*, 10:177–183.
- Reinsch, C. H. (1971). Smoothing by spline functions. ii. *Numerische Mathematik*, 16(5):451–454.
- Reiss, P. T. and Ogden, T. R. (2009). Smoothing parameter selection for a class of semiparametric linear models. *Journal of the Royal Statistical Society: Series B*, 71(2):505–523.
- Rue, H. and Held, L. (2005). *Gaussian Markov Random Fields: Theory and Applications*, volume 104 of *Monographs on Statistics and Applied Probability*. Chapman and Hall.
- Ruppert, D., Wand, M. P., and Carroll, R. J. (2003). *Semiparametric Regression*. Cambridge University Press.
- Shaddick, G. and Zidek, J. (2014). A case study in preferential sampling: Long term monitoring of air pollution in the UK. *Spatial Statistics*, 9.
- Silverman, B. (1985). Some aspects of the spline smoothing approach to non-parametric regression curve fitting. *Journal of the Royal Statistical Society: series B*, 47(1):1–53.
- Smith, M., Wong, C.-M., and Kohn, R. (1998). Additive nonparametric regression with autocorrelated errors. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 60(2):311–331.
- Spiegelhalter, D. J., Thomas, A., Best, N. G., and Gilks, W. R. (1996). *BUGS: Bayesian inference Using Gibbs Sampling, Version 0.5, (version ii)*. MRC Biostatistics Unit, Cambridge.
- Sven, H. and Craig, L. (2008). Updating the QR factorization and the least squares problem. Technical report, Manchester Institute for Mathematical Sciences, School of Mathematics, The University of Manchester. [http://eprints.ma.man.ac.uk/1192/1/grupdating\\_12nov08.pdf](http://eprints.ma.man.ac.uk/1192/1/grupdating_12nov08.pdf).
- Terrain-50 (2015). *OS Terrain 50: user guide and technical specification*. Ordnance Survey, Adanac Drive, Southampton, SO16 0AS.
- Tuerlinckx, F., Rijmen, F., Verbeke, G., and Boeck, P. (2010). Statistical inference in generalized linear mixed models: A review. *British Journal of Mathematical and Statistical Psychology*, 59(2):225–255.
- van de Geijn, R. and Quintana-Ortí, E. (2008). *The Science of programming matrix computations*. Lulu Press, Inc.
- Van Zee, F. G., Smith, T., Igual, F. D., Smelyanskiy, M., Zhang, X., Kistler, M., Austel, V., Gunnels, J., Low, T. M., Marker, B., Killough, L., and van de Geijn, R. A. (2016). The BLIS framework: Experiments in portability. *ACM Transactions on Mathematical Software*, 42(2):12:1–12:19.
- Van Zee, F. G. and van de Geijn, R. A. (2015). BLIS: A framework for rapidly instantiating BLAS functionality. *ACM Transactions on Mathematical Software*, 41(3):14:1–14:33.
- Wahba, G. (1980). *Approximation Theory III*, chapter Spline bases, regularization, and generalized cross validation for solving approximation problems with large quantities of noisy data, pages 905–912. Academic Press.
- Wahba, G. (1981). Spline interpolation and smoothing on the sphere. *SIAM Journal on Scientific and Statistical Computing*, 2(1):5–16.
- Wahba, G. (1983). Bayesian confidence intervals for the cross validated smoothing spline. *Journal of the Royal Statistical Society: series B*, 45:133–150.
- Wahba, G. (1990). *Spline Models for Observational Data*. Society for Industrial and Applied Mathematics.
- Wand, M. P. and Jones, M. C. (1994). *Kernel smoothing*, volume 60 of *Monographs on Statistics and Applied Probability*. Chapman and Hall/CRC.
- Wang, X., Shen, J., and Ruppert, D. (2011). On the asymptotics of penalized spline smoothing. *Electronic Journal of Statistics*, 5:1–17.



- Wang, Y. (1998). Smoothing spline models with correlated random errors. *Journal of the American Statistical Association*, 93(441):341–348.
- Watson, G. S. (1964). Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics, Series A (1961-2002)*, 26(4):359–372.
- Welford, B. P. (1962). Note on a method for calculating corrected sums of squares and products. *Technometrics*, 4(3):419–420.
- Whaley, R. C., Petitet, A. P., and Dongarra, J. J. (2000). Automated empirical optimization of software and the ATLAS project. Technical Report 147, LAPACK Working Notes.
- Wood, S. (2016). Just another gibbs additive modeler: Interfacing jags and mgcv. *Journal of Statistical Software, Articles*, 75(7):1–15.
- Wood, S., Li, Z., Shaddick, G., and Augustin, N. (2017). Generalized additive models for gigadata: modelling the UK Black Smoke Network daily data. *Journal of the American Statistical Association*, 112(519):1199–1210.
- Wood, S. N. (2003). Thin plate regression splines. *Journal of the Royal Statistical society: Series B*, 65(1):95–114.
- Wood, S. N. (2004). Stable and efficient multiple smoothing parameter estimation for generalized additive models. *Journal of the American Statistical Association*, 99(467):637–686.
- Wood, S. N. (2006). Low-rank scale-invariant tensor product smooths for generalized additive mixed models. *Biometrics*, 62(4):1025–1036.
- Wood, S. N. (2008). Fast stable direct fitting and smoothness selection for generalized additive models. *Journal of the Royal Statistical society: Series B*, 70(3):495–518.
- Wood, S. N. (2011). Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *Journal of the Royal Statistical society: Series B*, 73(1):3–36.
- Wood, S. N. (2017). P-splines with derivative based penalties and tensor product smoothing of unevenly distributed data. *Statistics and Computing*, 27(4):985–989.
- Wood, S. N., Bravington, M. V., and L., H. S. (2008). Soap film smoothing. *Journal of the Royal Statistical society: Series B*, 70(5):931–955.
- Wood, S. N., Goude, Y., and Shaw, S. (2015). Generalized additive models for large data sets. *Journal of the Royal Statistical society: Series C*, 64(1):139–155.
- Wood, S. N., Scheipl, F., and Faraway, J. (2013). Straightforward intermediate rank tensor product smoothing in mixed models. *Statistics and Computing*, 23(3):341–360.
- Zhang, X., Wang, Q., and Zhang, Y. (2011). OpenBLAS: A high performance BLAS library on Loongson 3A CPU. *Journal of Software*, 22:208–216.
- Zuur, A., Ieno, E., and Saveliev, A. (2014). *A Beginner’s Guide to Generalised Additive Mixed Models with R*. Highland Statistics Limited.